

INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(51) International Patent Classification ⁶ : C12N		A2	(11) International Publication Number: WO 98/59034 (43) International Publication Date: 30 December 1998 (30.12.98)
<p>(21) International Application Number: PCT/US98/13041</p> <p>(22) International Filing Date: 23 June 1998 (23.06.98)</p> <p>(30) Priority Data: 60/050,667 24 June 1997 (24.06.97) US</p> <p>(71) Applicant (for all designated States except US): HUMAN GENOME SCIENCES, INC. [US/US]; 9410 Key West Avenue, Rockville, MD 20850 (US).</p> <p>(72) Inventor; and</p> <p>(75) Inventor/Applicant (for US only): FRASER, Claire, M. [US/US]; 11915 Glen Mill Road, Potomac, MD 20854 (US).</p> <p>(74) Agents: BROOKES, A., Anders et al.; Human Genome Sciences, Inc., 9410 Key West Avenue, Rockville, MD 20850 (US).</p>		<p>(81) Designated States: AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, CA, CH, CN, CU, CZ, DE, DK, EE, ES, FI, GB, GE, GH, GM, GW, HU, ID, IL, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MD, MG, MK, MN, MW, MX, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, UA, UG, US, UZ, VN, YU, ZW, ARIPO patent (GH, GM, KE, LS, MW, SD, SZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, ML, MR, NE, SN, TD, TG).</p> <p>Published Without international search report and to be republished upon receipt of that report.</p>	
<p>(54) Title: TREPONEMA PALLIDUM POLYNUCLEOTIDES AND SEQUENCES</p> <pre> graph TD BUS[BUS] --> Processor[Processor] BUS --> MainMemory[Main Memory] MainMemory --> SSD[Secondary Storage Devices 110] SSD --- HardDrive[Hard Drive] SSD --- RMSD[Removable Medium Storage Device] RMSD -.-> RSM[Removable Storage Medium] </pre>			
<p>(57) Abstract</p> <p>The present invention provides polynucleotide sequences of the genome of <i>T. pallidum</i>, polypeptide sequences encoded by the polynucleotide sequences, corresponding polynucleotides and polypeptides, vectors and hosts comprising the polynucleotides, and assays and other uses thereof. The present invention further provides polynucleotide and polypeptide sequence information stored on computer readable media, and computer-based systems and methods which facilitate its use.</p>			

FOR THE PURPOSES OF INFORMATION ONLY

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AL	Albania	ES	Spain	LS	Lesotho	SI	Slovenia
AM	Armenia	FI	Finland	LT	Lithuania	SK	Slovakia
AT	Austria	FR	France	LU	Luxembourg	SN	Senegal
AU	Australia	GA	Gabon	LV	Latvia	SZ	Swaziland
AZ	Azerbaijan	GB	United Kingdom	MC	Monaco	TD	Chad
BA	Bosnia and Herzegovina	GE	Georgia	MD	Republic of Moldova	TG	Togo
BB	Barbados	GH	Ghana	MG	Madagascar	TJ	Tajikistan
BE	Belgium	GN	Guinea	MK	The former Yugoslav Republic of Macedonia	TM	Turkmenistan
BF	Burkina Faso	GR	Greece	ML	Mali	TR	Turkey
BG	Bulgaria	HU	Hungary	MN	Mongolia	TT	Trinidad and Tobago
BJ	Benin	IE	Ireland	MR	Mauritania	UA	Ukraine
BR	Brazil	IL	Israel	MW	Malawi	UG	Uganda
BY	Belarus	IS	Iceland	MX	Mexico	US	United States of America
CA	Canada	IT	Italy	NE	Niger	UZ	Uzbekistan
CF	Central African Republic	JP	Japan	NL	Netherlands	VN	Viet Nam
CG	Congo	KE	Kenya	NO	Norway	YU	Yugoslavia
CH	Switzerland	KG	Kyrgyzstan	NZ	New Zealand	ZW	Zimbabwe
CI	Côte d'Ivoire	KP	Democratic People's Republic of Korea	PL	Poland		
CM	Cameroon	KR	Republic of Korea	PT	Portugal		
CN	China	KZ	Kazakhstan	RO	Romania		
CU	Cuba	LC	Saint Lucia	RU	Russian Federation		
CZ	Czech Republic	LJ	Liechtenstein	SD	Sudan		
DE	Germany	LK	Sri Lanka	SE	Sweden		
DK	Denmark	LR	Liberia	SG	Singapore		
EE	Estonia						

Treponema pallidum Polynucleotides and Sequences**FIELD OF THE INVENTION**

The present invention relates to the field of molecular biology. In particular, it relates to, 5 among other things, nucleotide sequences of *Treponema pallidum*, contigs, ORFs, fragments, probes, primers and related polynucleotides thereof, peptides and polypeptides encoded by the sequences, and uses of the polynucleotides and sequences thereof, such as in fermentation, polypeptide production, assays and pharmaceutical development, among others.

10 BACKGROUND OF THE INVENTION

Spirochetes are a family of motile, unicellular, spiral-shaped bacteria which share a number of structural characteristics. Three genera of the spirochetes are pathogenic in humans:

(a) *Treponema*, which includes the pathogens that cause syphilis (*T. pallidum*), yaws (*T. pertenue*), and pinta (*T. carateum*); (b) *Borrelia*, which includes the pathogens that cause epidemic and endemic relapsing fever and Lyme disease; and (c) *Leptospira*, which includes a wide variety of small spirochetes that cause mild to serious systemic human illness (Koff, A. B. and Rosen, T. J. Am. Acad. Dermatol. 29:519-535 (1993)). In 1986, more than 27,000 cases 15 of early infectious syphilis were diagnosed in the United States alone. Such statistics indicate that infection with *T. pallidum* is the largest source of human disease resulting from the 20 spirochetes.

T. pallidum is morphologically indistinguishable from several other pathogenic spirochetes, but, in general, treponemes and other spirochetes, are easily identifiable when compared to other bacteria. A key morphological characteristic of *T. pallidum*, and other 25 spirochetes, is the presence of a central protoplasmic cylinder composed primarily of peptidoglycan and one or more adjacent axial fibrils (also designated periplasmic flagella or endoflagella; Charon, N. W., et al., Res. Microbiol. 143:597-603 (1992)). These structures provide a source of corkscrew-like motion to the treponemes. In aqueous media, treponemes move in an apparently random fashion and, unlike the majority of motile bacteria, continue to 30 move in a more viscous medium. In tissues, treponemes are highly moldable to intercellular spaces; a characteristic which is thought to be mediated by the interactions of bacterial adhesins and cellular fibronectins.

Syphilis is the primary clinical manifestation of infection with *T. pallidum*. The clinical manifestations of syphilis can resemble many diseases. Syphilis is typically transmitted by 35 sexual contact, but can also be transmitted transplacentally. The infecting organism multiplies at the site of infection within 10 to 60 days postinfection and results in a primary ulcer-like lesion termed a chancre. A small number of organisms move from the primary lesion to the regional lymph nodes and establish small infectious centers termed satellite buboes. Organisms from

these locations enter the blood stream and result in a systemic infection (Goens, J. L., *et al.*, *Am. Fam. Physician* 50:1013-1020 (1994)).

The secondary stage of syphilis manifests itself as a widespread skin rash and begins between two and twelve weeks following the primary infection. During this stage, the infected individual often experiences a low grade fever coupled with swollen lymph nodes. Also during this period, lesions of various degrees of severity may develop in a number of physical locations including bone, liver, kidney, central nervous system (CNS), and other organs (Veeravahu, M. *Arch. Intern. Med.* 145:132-134 (1985)). Such secondary infections are highly infectious, but will, in time, subside spontaneously.

A third stage of syphilis occurs in approximately 30% of infected, but not treated, individuals. The third stage occurs several years following the first and second stages. The lesions which characterize the third stage of infection are minor in terms of the number of organisms, but may be severe in terms of tissue damage. Such lesions may result in necrosis, scar formation, general paresis, damage to aortic valves, permanent blindness, and other extensive tissue damage, all probably related to a delayed type hypersensitivity reaction by the host to the *T. pallidum* organisms (Scheck, D. N. and Hook, E. W. *3rd Infect. Dis. Clin. North Am.* 8:769-795 (1994)).

A further, and increasingly common, complication of syphilis infection is coinfection with the human immunodeficiency virus (HIV). In fact, a recent study indicates that ulcerous genital diseases such as those exhibited during the primary stages of infection with syphilis may facilitate the transmission of HIV (Rufli, T. *Dermatologica* 179:113-117 (1989)). In addition, it is clear that the CNS is regularly involved in the early stages of syphilis. In the timespan between the introduction of penecillin and other antibiotics and the spread of HIV, early neurosyphilis was an exceptionally uncommon development. However, since the standard antibiotic dosage used to treat syphilis is not exceptionally high and since a successful treatment requires an adequate host immune response, individuals infected with HIV often exhibit a highly increased occurrence of many neurosyphilis-related sequelae including asymptomatic neurosyphilis, syphilitic meinigitis, cranial nerve abnormalities, or cerebrovascular problems (Musher, D. M., *et al.*, *Ann. Intern. Med.* 113:872-881 (1990)).

T. pallidum has a remarkable ability to evade both the humoral and cellular components of the immune system. It was originally thought that the ability of *T. pallidum* to evade the immune system of the host organism was due to the presence of an outer coat of mucopolysaccharides. However, recent evidence suggests it is more likely that *T. pallidum* make use of the organization of the relative immunogenicity of its complement of outer membrane proteins to evade the immune system (Radolf, J. D. *Mol. Microbiol.* 16:1067-1073 (1995)). Unlike most other bacterial outer membranes characterized thus far, the *T. pallidum* outer membrane contains a scarcity of immunogenic transmembrane proteins (with regard to *T. pallidum*, these are termed "rare outer membrane proteins"). Among the highly immunogenic proteins of treponemes are a number of lipoproteins anchored to the periplasmic leaflet of the cytoplasmic membrane. As a

result of their physical location, the lipoproteins may be less susceptible to typical immunologic surveillance (Norris, *J. Microbiol. Rev.* 57:750-779 (1993)). In addition to the periplasmic lipoproteins, *T. pallidum* also secretes a number of small, but immunogenic proteins which may induce an immune response (Hindersson, P. et al., *Res. Microbiol.* 143:629-639 (1992)).

- 5 It is clear that the etiology of diseases mediated or exacerbated by *T. pallidum* genes, and that characterizing the genes and their patterns of expression would add dramatically to our understanding of the organism and its host interactions. Knowledge of *T. pallidum* genes and genomic organization would dramatically improve understanding of disease etiology and lead to improved and new ways of preventing, ameliorating, arresting and reversing diseases.
- 10 Moreover, characterized genes and genomic fragments of *T. pallidum* would provide reagents for, among other things, detecting, characterizing and controlling *T. pallidum* infections. There is a need therefore to characterize the genome of *T. pallidum* and for polynucleotides and sequences of this organism.

15 **SUMMARY OF THE INVENTION**

The present invention is based on the sequencing of fragments of the *T. pallidum* genome. The primary nucleotide sequences which were generated are provided in SEQ ID NOS:1-744.

- 20 The present invention provides the nucleotide sequence of several thousand contigs of the *T. pallidum* genome, which are listed in tables below and set out in the Sequence Listing submitted herewith, and representative fragments thereof, in a form which can be readily used, analyzed, and interpreted by a skilled artisan. In one embodiment, the present invention is provided as contiguous strings of primary sequence information corresponding to the nucleotide sequences depicted in SEQ ID NOS: 1-744.

- 25 The present invention further provides nucleotide sequences which are at least 95% identical to the nucleotide sequences of SEQ ID NOS: 1-744.

- 30 The nucleotide sequence of SEQ ID NOS: 1-744 , a representative fragment thereof, or a nucleotide sequence which is at least 95% identical to the nucleotide sequence of SEQ ID NOS: 1-744 may be provided in a variety of mediums to facilitate its use. In one application of this embodiment, the sequences of the present invention are recorded on computer readable media. Such media includes, but is not limited to: magnetic storage media, such as floppy discs, hard disc storage medium, and magnetic tape; optical storage media such as CD-ROM; electrical storage media such as RAM and ROM; and hybrids of these categories such as magnetic/optical storage media.

- 35 The present invention further provides systems, particularly computer-based systems which contain the sequence information herein described stored in a data storage means. Such systems are designed to identify commercially important fragments of the *T. pallidum* genome.

Another embodiment of the present invention is directed to fragments of the *T. pallidum* genome having particular structural or functional attributes. Such fragments of the *T. pallidum*

genome of the present invention include, but are not limited to, fragments which encode peptides, hereinafter referred to as open reading frames or ORFs, fragments which modulate the expression of an operably linked ORF, hereinafter referred to as expression modulating fragments or EMFs, and fragments which can be used to diagnose the presence of *T. pallidum* 5 in a sample, hereinafter referred to as diagnostic fragments or DFs.

Each of the ORFs in fragments of the *T. pallidum* genome disclosed in Tables 1, 2 and 3, and the EMFs found 5' to the ORFs, can be used in numerous ways as polynucleotide reagents. For instance, the sequences can be used as diagnostic probes or amplification primers for detecting or determining the presence of a specific microbe in a sample, to selectively control 10 gene expression in a host and in the production of polypeptides, such as polypeptides encoded by ORFs of the present invention, particular those polypeptides that have a pharmacological activity.

The present invention further includes recombinant constructs comprising one or more fragments of the *T. pallidum* genome of the present invention. The recombinant constructs of the 15 present invention comprise vectors, such as a plasmid or viral vector, into which a fragment of the *T. pallidum* has been inserted.

The present invention further provides host cells containing any of the isolated fragments of the *T. pallidum* genome of the present invention. The host cells can be a higher eukaryotic host cell, such as a mammalian cell, a lower eukaryotic cell, such as a yeast cell, or a procaryotic cell such as a bacterial cell.

20 The present invention is further directed to isolated polypeptides and proteins encoded by ORFs of the present invention. A variety of methods, well known to those of skill in the art, routinely may be utilized to obtain any of the polypeptides and proteins of the present invention. For instance, polypeptides and proteins of the present invention having relatively short, simple amino acid sequences readily can be synthesized using commercially available automated peptide 25 synthesizers. Polypeptides and proteins of the present invention also may be purified from bacterial cells which naturally produce the protein. Yet another alternative is to purify polypeptide and proteins of the present invention from cells which have been altered to express them.

30 The invention further provides methods of obtaining homologs of the fragments of the *T. pallidum* genome of the present invention and homologs of the proteins encoded by the ORFs of the present invention. Specifically, by using the nucleotide and amino acid sequences disclosed herein as a probe or as primers, and techniques such as PCR cloning and colony/plaque hybridization, one skilled in the art can obtain homologs.

35 The invention further provides antibodies which selectively bind polypeptides and proteins of the present invention. Such antibodies include both monoclonal and polyclonal antibodies.

The invention further provides hybridomas which produce the above-described antibodies. A hybridoma is an immortalized cell line which is capable of secreting a specific monoclonal antibody.

The present invention further provides methods of identifying test samples derived from cells which express one of the ORFs of the present invention, or a homolog thereof. Such methods comprise incubating a test sample with one or more of the antibodies of the present invention, or one or more of the DFs of the present invention, under conditions which allow a skilled artisan to determine if the sample contains the ORF or product produced therefrom.

In another embodiment of the present invention, kits are provided which contain the necessary reagents to carry out the above-described assays.

Specifically, the invention provides a compartmentalized kit to receive, in close confinement, one or more containers which comprises: (a) a first container comprising one of the antibodies, or one of the DFs of the present invention; and (b) one or more other containers comprising one or more of the following: wash reagents, reagents capable of detecting presence of bound antibodies or hybridized DFs.

Using the isolated proteins of the present invention, the present invention further provides methods of obtaining and identifying agents capable of binding to a polypeptide or protein encoded by one of the ORFs of the present invention. Specifically, such agents include, as further described below, antibodies, peptides, carbohydrates, pharmaceutical agents and the like. Such methods comprise steps of: (a) contacting an agent with an isolated protein encoded by one of the ORFs of the present invention; and (b) determining whether the agent binds to said protein.

The present genomic sequences of *T. pallidum* will be of great value to all laboratories working with this organism and for a variety of commercial purposes. Many fragments of the *T. pallidum* genome will be immediately identified by similarity searches against GenBank or protein databases and will be of immediate value to *T. pallidum* researchers and for immediate commercial value for the production of proteins or to control gene expression.

The methodology and technology for elucidating extensive genomic sequences of bacterial and other genomes has and will greatly enhance the ability to analyze and understand chromosomal organization. In particular, sequenced contigs and genomes will provide the models for developing tools for the analysis of chromosome structure and function, including the ability to identify genes within large segments of genomic DNA, the structure, position, and spacing of regulatory elements, the identification of genes with potential industrial applications, and the ability to do comparative genomic and molecular phylogeny.

DESCRIPTION OF THE FIGURES

FIGURE 1 is a block diagram of a computer system (102) that can be used to implement computer-based systems of present invention.

FIGURE 2 is a schematic diagram depicting the data flow and computer programs used to collect, assemble, edit and annotate the contigs of the *T. pallidum* genome of the present invention: Both Macintosh and Unix platforms are used to handle the AB 373 and 377 sequence

data files, largely as described in Kerlavage *et al.*, *Proceedings of the Twenty-Sixth Annual Hawaii International Conference on System Sciences*, 585, IEEE Computer Society Press, Washington D.C. (1993). Factura (AB) is a Macintosh program designed for automatic vector sequence removal and end-trimming of sequence files. The program Loadis runs on a Macintosh platform and parses the feature data extracted from the sequence files by Factura to the Unix based *T. pallidum* relational database. Assembly of contigs (and whole genome sequences) is accomplished by retrieving a specific set of sequence files and their associated features using Extrseq, a Unix utility for retrieving sequences from an SQL database. The resulting sequence file is processed to trim portions of the sequences with a high rate ambiguous nucleotides. The sequence files were assembled using TIGR Assembler, an assembly engine designed at The Institute for Genomic Research (TIGR) for rapid and accurate assembly of thousands of sequence fragments. The collection of contigs generated by the assembly step is loaded into the database with the lassie program. Identification of open reading frames (ORFs) is accomplished by processing contigs with zorf. The ORFs are searched against *T. pallidum* sequences from GenBank and against all protein sequences using the BLASTN and BLASTP programs (using default parameters), described in Altschul *et al.*, *J. Mol. Biol.* 215: 403-410 (1990). Results of the ORF determination and similarity searching steps were loaded into the database. As described below, some results of the determination and the searches are set out in Tables 1-3.

20 **DETAILED DESCRIPTION OF ILLUSTRATIVE EMBODIMENTS**

The present invention is based on the sequencing of fragments of the *T. pallidum* genome and analysis of the sequences. The primary nucleotide sequences generated by sequencing the fragments are provided in SEQ ID NOS: 1-744. As used herein, the "primary sequence" refers to the nucleotide sequence represented by the IUPAC nomenclature system.).

25 In addition to the aforementioned *T. pallidum* polynucleotide and polynucleotide sequences, the present invention provides the nucleotide sequences of SEQ ID NOS: 1-744, ORF IDs and ORFs within, or representative fragments thereof, in a form which can be readily used, analyzed, and interpreted by a skilled artisan.

As used herein, a "representative fragment of the nucleotide sequence depicted in SEQ ID 30 NOS:1-744" refers to any portion of the SEQ ID NOS: 1-744 which is not presently represented within a publicly available database. Preferred representative fragments of the present invention are *T. pallidum* open reading frames (ORFs), expression modulating fragment (EMFs) and fragments which can be used to diagnose the presence of *T. pallidum* in sample (DFs). A non-limiting identification of preferred representative fragments is provided in Tables 1-3. As 35 discussed in detail below, the information provided in SEQ ID NOS:1-744 and in Tables 1-3 together with routine cloning, synthesis, sequencing and assay methods will enable those skilled in the art to clone and sequence all "representative fragments" of interest, including open reading frames encoding a large variety of *T. pallidum* proteins.

The present invention is further directed to nucleic acid molecules encoding portions or fragments of the nucleotide sequences described herein. Fragments include portions of the nucleotide sequences of SEQ ID NOS:1-744, at least 10 contiguous nucleotides in length selected from any two integers, one of which representing a 5' nucleotide position and a second of which 5 representing a 3' nucleotide position, where the first nucleotide for each nucleotide sequence in SEQ ID NOS:1-744 is position 1. That is, every combination of a 5' and 3' nucleotide position that a fragment at least 10 contiguous nucleotides in length could occupy is included in the invention. At least means a fragment may be 10 contiguous nucleotide bases in length or any integer between 10 and the length of an entire nucleotide sequence of SEQ ID NOS:1-744 minus 10 1. Therefore, included in the invention are contiguous fragments specified by any 5' and 3' nucleotide base positions of a nucleotide sequences of SEQ ID NOS:1-744 wherein the contiguous fragment is any integer between 10 and the length of an entire nucleotide sequence minus 1.

Further, the invention includes polynucleotides comprising fragments specified by size, 15 in nucleotides, rather than by nucleotide positions. The invention includes any fragment size, in contiguous nucleotides, selected from integers between 10 and the length of an entire ORF ID, ORF, or SEQ ID NO:, minus 1. Preferred sizes of contiguous nucleotide fragments include 20 nucleotides, 30 nucleotides, 40 nucleotides, 50 nucleotides. Other preferred sizes of contiguous nucleotide fragments, which may be useful as diagnostic probes and primers, include fragments 20 50-300 nucleotides in length which include, as discussed above, fragment sizes representing each integer between 50-300. Larger fragments are also useful according to the present invention corresponding to most, if not all, of the nucleotide sequences shown in Tables 1-3 (ORF IDs) and SEQ ID NOS:1-744. The preferred sizes are, of course, meant to exemplify not limit the present invention as all size fragments, representing any integer between 10 and the length of an 25 entire nucleotide sequence minus 1, of each ORF ID, ORF, and SEQ ID NO:, are included in the invention.

The present invention also provides for the exclusion of any fragment, specified by 5' and 3' base positions or by size in nucleotide bases as described above for any ORF ID or SEQ 30 ID NOS:1-744. Any number of fragments of nucleotide sequences in ORF IDs or SEQ ID NOS:1-744, specified by 5' and 3' base positions or by size in nucleotides, as described above, may be excluded from the present invention.

While the presently disclosed sequences of SEQ ID NOS: 1-744 are highly accurate, sequencing techniques are not perfect and, in relatively rare instances, further investigation of a fragment or sequence of the invention may reveal a nucleotide sequence error present in a 35 nucleotide sequence disclosed in SEQ ID NOS: 1-744. However, once the present invention is made available (*i.e.*, once the information in SEQ ID NOS: 1-744 and Tables 1-3 has been made available), resolving a rare sequencing error in SEQ ID NOS: 1-744 will be well within the skill of the art. The present disclosure makes available sufficient sequence information to allow any of the described contigs or portions thereof to be obtained readily by straightforward application of

routine techniques. Further sequencing of such polynucleotide may proceed in like manner using manual and automated sequencing methods which are employed ubiquitous in the art. Nucleotide sequence editing software is publicly available. For example, Applied Biosystem's (AB) AutoAssembler can be used as an aid during visual inspection of nucleotide sequences. By 5 employing such routine techniques potential errors readily may be identified and the correct sequence then may be ascertained by targeting further sequencing effort, also of a routine nature, to the region containing the potential error.

Even if all of the very rare sequencing errors in SEQ ID NOS: 1-744 were corrected, the resulting nucleotide sequences would still be at least 95% identical, nearly all would be at least 10 99% identical, and the great majority would be at least 99.9% identical to the nucleotide sequences of SEQ ID NOS: 1-7441-744.

As discussed elsewhere herein, polynucleotides of the present invention readily may be obtained by routine application of well known and standard procedures for cloning and sequencing DNA. Detailed methods for obtaining libraries and for sequencing are provided 15 below, for instance. A wide variety of *T. pallidum* strains can be used to prepare *T. pallidum* genomic DNA for cloning and for obtaining polynucleotides of the present invention which are known in the art.

The nucleotide sequences of the genomes from different strains of *T. pallidum* differ somewhat. However, the nucleotide sequences of the genomes of all *T. pallidum* strains will be 20 at least 95% identical, in corresponding part, to the nucleotide sequences provided in SEQ ID NOS: 1-744 and the ORF IDs and ORFs within. Nearly all will be at least 99% identical and the great majority will be 99.9% identical.

The present application is further directed to nucleic acid molecules at least 90%, 95%, 96%, 97%, 98% or 99% identical to a nucleic acid sequence shown in SEQ ID NOS: 1-744, the 25 ORF IDs and ORFs within. The above nucleic acid sequences are included irrespective of whether they encode a polypeptide having *T. pallidum* activity. This is because even where a particular nucleic acid molecule does not encode a polypeptide having *T. pallidum* activity, one of skill in the art would still know how to use the nucleic acid molecule, for instance, as a hybridization probe. Uses of the nucleic acid molecules of the present invention that do not 30 encode a polypeptide having *T. pallidum* activity include, *inter alia*, isolating an *T. pallidum* gene or allelic variants thereof from a DNA library, and detecting *T. pallidum* mRNA expression samples, environmental samples, suspected of containing *T. pallidum* by Northern Blot, PCR, or similar analysis.

Preferred, are nucleic acid molecules having sequences at least 90%, 95%, 96%, 97%, 35 98% or 99% identical to the nucleic acid sequence shown in SEQ ID NOS: 1-744, the ORF IDs, and the ORF within each ORF ID, which do, in fact, encode a polypeptide having *T. pallidum* protein activity. By "a polypeptide having *T. pallidum* activity" is intended polypeptides exhibiting activity similar, but not necessarily identical, to an activity of the *T. pallidum* protein of the invention, as measured in a particular biological assay suitable for measuring activity of the

specified protein.

Due to the degeneracy of the genetic code, one of ordinary skill in the art will immediately recognize that a large number of the nucleic acid molecules having a sequence at least 90%, 95%, 96%, 97%, 98%, or 99% identical to the nucleic acid sequences shown in SEQ ID NOS: 1-744,

5 the ORF IDs, and the ORF within each ORF ID, will encode a polypeptide having *T. pallidum* protein activity. In fact, since degenerate variants of these nucleotide sequences all encode the same polypeptide, this will be clear to the skilled artisan even without performing the above described comparison assay. It will be further recognized in the art that, for such nucleic acid molecules that are not degenerate variants, a reasonable number will also encode a polypeptide
10 having *T. pallidum* protein activity. This is because the skilled artisan is fully aware of amino acid substitutions that are either less likely or not likely to significantly effect protein function (e.g., replacing one aliphatic amino acid with a second aliphatic amino acid), as further described below.

The biological activity or function of the polypeptides of the present invention are
15 expected to be similar or identical to polypeptides from other bacteria that share a high degree of structural identity/similarity. Table 1-3 lists accession numbers and descriptions for the closest matching sequences of polypeptides available through Genbank. It is therefore expected that the biological activity or function of the polypeptides of the present invention will be similar or identical to those polypeptides from other bacterial genera, species, or strains listed in Table 1-
20 3.

By a polynucleotide having a nucleotide sequence at least, for example, 95% "identical" to a reference nucleotide sequence of the present invention, it is intended that the nucleotide sequence of the polynucleotide is identical to the reference sequence except that the polynucleotide sequence may include up to five point mutations per each 100 nucleotides of the
25 reference nucleotide sequence encoding the *T. pallidum* polypeptide. In other words, to obtain a polynucleotide having a nucleotide sequence at least 95% identical to a reference nucleotide sequence, up to 5% of the nucleotides in the reference sequence may be deleted, inserted, or substituted with another nucleotide. The query sequence may be an entire sequence shown in SEQ ID NOS: 1-744, the ORF IDs, or the ORF within each ORF ID, or any fragment specified
30 as described herein.

As a practical matter, whether any particular nucleic acid molecule or polypeptide is at least 90%, 95%, 96%, 97%, 98% or 99% identical to a nucleotide sequence of the presence invention can be determined conventionally using known computer programs. A preferred method for determining the best overall match between a query sequence (a sequence of the
35 present invention) and a subject sequence, also referred to as a global sequence alignment, can be determined using the FASTDB computer program based on the algorithm of Brutlag et al. See Brutlag et al. (1990) Comp. App. Biosci. 6:237-245. In a sequence alignment the query and subject sequences are both DNA sequences. An RNA sequence can be compared by first converting U's to T's. The result of said global sequence alignment is in percent identity.

Preferred parameters used in a FASTDB alignment of DNA sequences to calculate percent identity are: Matrix=Unitary, k-tuple=4, Mismatch Penalty=1, Joining Penalty=30, Randomization Group Length=0, Cutoff Score=1, Gap Penalty=5, Gap Size Penalty 0.05, Window Size=500 or the lenght of the subject nucleotide sequence, whichever is shorter.

- 5 If the subject sequence is shorter than the query sequence because of 5' or 3' deletions, not because of internal deletions, a manual correction must be made to the results. This is because the FASTDB program does not account for 5' and 3' truncations of the subject sequence when calculating percent identity. For subject sequences truncated at the 5' or 3' ends, relative to the query sequence, the percent identity is corrected by calculating the number of bases of the
- 10 query sequence that are 5' and 3' of the subject sequence, which are not matched/aligned, as a percent of the total bases of the query sequence. Whether a nucleotide is matched/aligned is determined by results of the FASTDB sequence alignment. This percentage is then subtracted from the percent identity, calculated by the above FASTDB program using the specified parameters, to arrive at a final percent identity score. This corrected score is what is used for the
- 15 purposes of the present invention. Only nucleotides outside the 5' and 3' nucleotides of the subject sequence, as displayed by the FASTDB alignment, which are not matched/aligned with the query sequence, are calculated for the purposes of manually adjusting the percent identity score.
- For example, a 90 nucleotide subject sequence is aligned to a 100 nucleotide query sequence to
- 20 determine percent identity. The deletions occur at the 5' end of the subject sequence and therefore, the FASTDB alignment does not show a matched/alignment of the first 10 nucleotides at 5' end. The 10 unpaired nucleotides represent 10% of the sequence (number of nucleotides at the 5' and 3' ends not matched/total number of nucleotides in the query sequence) so 10% is subtracted from the percent identity score calculated by the FASTDB program. If the remaining
- 25 90 nucleotides were perfectly matched the final percent identity would be 90%. In another example, a 90 nucleotide subject sequence is compared with a 100 nucleotide query sequence. This time the deletions are internal deletions so that there are no nucleotides on the 5' or 3' of the subject sequence which are not matched/aligned with the query. In this case the percent identity calculated by FASTDB is not manually corrected. Once again, only nucleotides 5' and 3' of the
- 30 subject sequence which are not matched/aligned with the query sequence are manually corrected for. No other manual corrections are to be made for the purposes of the present invention.

COMPUTER RELATED EMBODIMENTS

The nucleotide sequences provided in SEQ ID NOS: 1-744, including ORF IDs and

35 corresponding ORFs, a representative fragment thereof, or a nucleotide sequence at least 95%, preferably at least 99% and most preferably at least 99.9% identical to said polynucleotide sequences may be "provided" in a variety of mediums to facilitate use thereof. As used herein, "provided" refers to a manufacture, other than an isolated nucleic acid molecule, which contains a nucleotide sequence of the present invention. Such a manufacture provides a large portion of the

T. pallidum genome and parts thereof (e.g., a *T. pallidum* open reading frame (ORF)) in a form which allows a skilled artisan to examine the manufacture using means not directly applicable to examining the *T. pallidum* genome or a subset thereof as it exists in nature or in purified form.

In one application of this embodiment, a nucleotide sequence of the present invention can

- 5 be recorded on computer readable media. As used herein, "computer readable media" refers to any medium which can be read and accessed directly by a computer. Such media include, but are not limited to: magnetic storage media, such as floppy discs, hard disc storage medium, and magnetic tape; optical storage media such as CD- ROM; electrical storage media such as RAM and ROM; and hybrids of these categories, such as magnetic/optical storage media. A skilled
- 10 artisan can readily appreciate how any of the presently known computer readable mediums can be used to create a manufacture comprising computer readable medium having recorded thereon a nucleotide sequence of the present invention. Likewise, it will be clear to those of skill how additional computer readable media that may be developed also can be used to create analogous manufactures having recorded thereon a nucleotide sequence of the present invention.

- 15 As used herein, "recorded" refers to a process for storing information on computer readable medium. A skilled artisan can readily adopt any of the presently known methods for recording information on computer readable medium to generate manufactures comprising the nucleotide sequence information of the present invention.

- 20 A variety of data storage structures are available to a skilled artisan for creating a computer readable medium having recorded thereon a nucleotide sequence of the present invention. The choice of the data storage structure will generally be based on the means chosen to access the stored information. In addition, a variety of data processor programs and formats can be used to store the nucleotide sequence information of the present invention on computer readable medium. The sequence information can be represented in a word processing text file, formatted in commercially-available software such as WordPerfect and MicroSoft Word, or represented in the form of an ASCII file, stored in a database application, such as DB2, Sybase, Oracle, or the like. A skilled artisan can readily adapt any number of data-processor structuring formats (e.g., text file or database) in order to obtain computer readable medium having recorded thereon the nucleotide sequence information of the present invention.
- 25

- 30 Computer software is publicly available which allows a skilled artisan to access sequence information provided in a computer readable medium. Thus, by providing in computer readable form the nucleotide sequences of SEQ ID NOS: 1-744, including ORF IDs and corresponding ORFs, a representative fragment thereof, or a nucleotide sequence at least 95%, preferably at least 99% and most preferably at least 99.9% identical to said polynucleotide sequences, the present invention enables the skilled artisan routinely to access the provided sequence information for a wide variety of purposes.
- 35

The examples which follow demonstrate how software which implements the BLAST (Altschul *et al.*, *J. Mol. Biol.* 215:403-410 (1990)) and BLAZE (Brutlag *et al.*, *Comp. Chem.* 17:203-207 (1993)) search algorithms on a Sybase system was used to identify open reading

frames (ORFs) within the *T. pallidum* genome which contain homology to ORFs or proteins from both *T. pallidum* and from other organisms. Among the ORFs discussed herein are protein encoding fragments of the *T. pallidum* genome useful in producing commercially important proteins, such as enzymes used in fermentation reactions and in the production of commercially useful metabolites.

The present invention further provides systems, particularly computer-based systems, which contain the sequence information described herein. Such systems are designed to identify, among other things, commercially important fragments of the *T. pallidum* genome.

As used herein, "a computer-based system" refers to the hardware means, software means, and data storage means used to analyze the nucleotide sequence information of the present invention. The minimum hardware means of the computer-based systems of the present invention comprises a central processing unit (CPU), input means, output means, and data storage means. A skilled artisan can readily appreciate that any one of the currently available computer-based system are suitable for use in the present invention.

As stated above, the computer-based systems of the present invention comprise a data storage means having stored therein a nucleotide sequence of the present invention and the necessary hardware means and software means for supporting and implementing a search means.

As used herein, "data storage means" refers to memory which can store nucleotide sequence information of the present invention, or a memory access means which can access manufactures having recorded thereon the nucleotide sequence information of the present invention.

As used herein, "search means" refers to one or more programs which are implemented on the computer-based system to compare a target sequence or target structural motif with the sequence information stored within the data storage means. Search means are used to identify fragments or regions of the present genomic sequences which match a particular target sequence or target motif. A variety of known algorithms are disclosed publicly and a variety of commercially available software for conducting search means are and can be used in the computer-based systems of the present invention. Examples of such software includes, but is not limited to, MacPattern (EMBL), BLASTN and BLASTX (NCBIA). A skilled artisan can readily recognize that any one of the available algorithms or implementing software packages for conducting homology searches can be adapted for use in the present computer-based systems.

As used herein, a "target sequence" can be any DNA or amino acid sequence of six or more nucleotides or two or more amino acids. A skilled artisan can readily recognize that the longer a target sequence is, the less likely a target sequence will be present as a random occurrence in the database. The most preferred sequence length of a target sequence is from about 10 to 100 amino acids or from about 30 to 300 nucleotide residues. However, it is well recognized that searches for commercially important fragments, such as sequence fragments involved in gene expression and protein processing, may be of shorter length.

As used herein, "a target structural motif," or "target motif," refers to any rationally selected sequence or combination of sequences in which the sequence(s) are chosen based on a three-dimensional configuration which is formed upon the folding of the target motif. There are a variety of target motifs known in the art. Protein target motifs include, but are not limited to, enzymic active sites and signal sequences. Nucleic acid target motifs include, but are not limited to, promoter sequences, hairpin structures and inducible expression elements (protein binding sequences).

A variety of structural formats for the input and output means can be used to input and output the information in the computer-based systems of the present invention. A preferred format for an output means ranks fragments of the *T. pallidum* genomic sequences possessing varying degrees of homology to the target sequence or target motif. Such presentation provides a skilled artisan with a ranking of sequences which contain various amounts of the target sequence or target motif and identifies the degree of homology contained in the identified fragment.

A variety of comparing means can be used to compare a target sequence or target motif with the data storage means to identify sequence fragments of the *T. pallidum* genome. In the present examples, implementing software which implement the BLAST and BLAZE algorithms, described in Altschul *et al.*, *J. Mol. Biol.* 215: 403-410 (1990), is used to identify open reading frames within the *T. pallidum* genome. A skilled artisan can readily recognize that any one of the publicly available homology search programs can be used as the search means for the computer-based systems of the present invention. Of course, suitable proprietary systems that may be known to those of skill also may be employed in this regard.

Figure 1 provides a block diagram of a computer system illustrative of embodiments of this aspect of present invention. The computer system 102 includes a processor 106 connected to a bus 104. Also connected to the bus 104 are a main memory 108 (preferably implemented as random access memory, RAM) and a variety of secondary storage devices 110, such as a hard drive 112 and a removable medium storage device 114. The removable medium storage device 114 may represent, for example, a floppy disk drive, a CD-ROM drive, a magnetic tape drive, etc. A removable storage medium 116 (such as a floppy disk, a compact disk, a magnetic tape, etc.) containing control logic and/or data recorded therein may be inserted into the removable medium storage device 114. The computer system 102 includes appropriate software for reading the control logic and/or the data from the removable medium storage device 114, once it is inserted into the removable medium storage device 114.

A nucleotide sequence of the present invention may be stored in a well known manner in the main memory 108, any of the secondary storage devices 110, and/or a removable storage medium 116. During execution, software for accessing and processing the genomic sequence (such as search tools, comparing tools, etc.) reside in main memory 108, in accordance with the requirements and operating parameters of the operating system, the hardware system and the software program or programs.

BIOCHEMICAL EMBODIMENTS

Other embodiments of the present invention are directed to isolated fragments of the *T. pallidum* genome. The fragments of the *T. pallidum* genome of the present invention include, but 5 are not limited to fragments which encode peptides, hereinafter open reading frames (ORFs), fragments which modulate the expression of an operably linked ORF, hereinafter expression modulating fragments (EMFs) and fragments which can be used to diagnose the presence of *T. pallidum* in a sample, hereinafter diagnostic fragments (DFs).

As used herein, an "isolated nucleic acid molecule" or an "isolated fragment of the *T. pallidum* genome" refers to a nucleic acid molecule possessing a specific nucleotide sequence 10 which has been subjected to purification means to reduce, from the composition, the number of compounds which are normally associated with the composition. Particularly, the term refers to the nucleic acid molecules having the sequences set out in SEQ ID NOS: 1-744, to representative fragments thereof as described above including ORF IDs and ORFs, to polynucleotides at least 15 95%, preferably at least 96%, 97%, 98%, or 99% and especially preferably at least 99.9% identical in sequence thereto, also as set out above.

A variety of purification means can be used to generate the isolated fragments of the present invention. These include, but are not limited to methods which separate constituents of a solution based on charge, solubility, or size.

20 In one embodiment, *T. pallidum* DNA can be enzymatically sheared to produce fragments of 15-20 kb in length. These fragments can then be used to generate a *T. pallidum* library by inserting them into lambda clones as described in the Examples below. Primers flanking, for example, an ORF, such as those enumerated in the ORF IDs of Tables 1-3, can then be generated using nucleotide sequence information provided in SEQ ID NOS: 1-744. Well known and 25 routine techniques of PCR cloning then can be used to isolate the ORF from the lambda DNA library or *T. pallidum* genomic DNA. Thus, given the availability of SEQ ID NOS:1-744, the information in Tables 1, 2 and 3, and the information that may be obtained readily by analysis of the sequences of SEQ ID NOS:1-744 using methods set out above, those of skill will be enabled by the present disclosure to isolate any ORF-containing or other nucleic acid fragment of the 30 present invention.

The isolated nucleic acid molecules of the present invention include, but are not limited to single stranded and double stranded DNA, and single stranded RNA. For purposes of numbering and reference to polynucleotide and polypeptide sequences the entire sequence of each sequence of SEQ ID NOS:1-744 is included with the first nucleotide being position 1. 35 Therefore, for reference purposes the numbering used in the present invention is that provided in the sequence listing for SEQ ID NOS:1-744.

As used herein, an open reading frame (ORF), means a series of nucleotide triplets coding for amino acid residues without any termination codons and is a sequence translatable into protein. Further, unless specified, the term "ORF" for each ORF ID is defined by the termination

codon at the 3' end and the 5' most methionine codon, at the 5' end, in frame with said 3' termination codon. Unless specified, the term "ORF" also refers to a particular polypeptide sequence defined by the ORF polynucleotide sequence, wherein the N-terminus is defined by the 5' most methionine codon in frame with the termination codon at the 3' end of the ORF ID and the C-terminus is defined by the last codon before the said 3' termination codon. As used herein, an ORF ID represents a sequence without any internal termination codons flanked by termination codons.

Tables 1, 2, and 3 list ORF IDs in the *T. pallidum* genomic contigs of the present invention that were identified as putative coding regions by the GeneMark software using organism-specific second-order Markov probability transition matrices. It will be appreciated that other criteria can be used, in accordance with well known analytical methods, such as those discussed herein, to generate more inclusive, more restrictive, or more selective lists.

Table 1 sets out ORF IDs in the *T. pallidum* contigs of the present invention that over a continuous region of at least 50 bases are 95% or more identical (by BLAST analysis) to a nucleotide sequence available through GenBank in June, 1997.

Table 2 sets out ORF IDs in the *T. pallidum* contigs of the present invention that are not in Table 1 and match, with a BLASTP probability score of 0.01 or less, a polypeptide sequence available through GenBank in July, 1996.

Table 3 sets out ORF IDs in the *T. pallidum* contigs of the present invention that do not match significantly, by BLASTP analysis, a polypeptide sequence available through GenBank in July, 1996.

In each table, the first and second columns identify the ORF ID by, respectively, contig number and ORF ID number within the contig; the third column indicates the first nucleotide of the ORF ID, counting from the 5' end of the contig strand; and the fourth column indicates the last nucleotide of the ORF ID, counting from the 5' end of the contig strand.

In Tables 1 and 2, column six, lists the Reference for the closest matching sequence available through GenBank. These reference numbers are the databases entry numbers commonly used by those of skill in the art, who will be familiar with their denominators. Descriptions of the nomenclature are available from the National Center for Biotechnology Information. Column seven in Tables 1 and 2 provides the gene name of the matching sequence; column eight provides the BLAST identity score from the comparison of the ORF and the homologous gene; and column nine indicates the length in nucleotides of the highest scoring segment pair identified by the BLAST identity analysis.

In Table 3, the last column, column six, indicates the length of each ORF ID in amino acid residues.

The concepts of percent identity and percent similarity of two polypeptide sequences is well understood in the art. For example, two polypeptides 10 amino acids in length which differ at three amino acid positions (e.g., at positions 1, 3 and 5) are said to have a percent identity of 70%. However, the same two polypeptides would be deemed to have a percent similarity of

80% if, for example at position 5, the amino acids moieties, although not identical, were "similar" (*i.e.*, possessed similar biochemical characteristics). Many programs for analysis of nucleotide or amino acid sequence similarity, such as FASTA and BLAST specifically list percent identity of a matching region as an output parameter. Thus, for instance, Tables 1 and 2 herein 5 enumerate the percent identity of the highest scoring segment pair in each ORF and its listed relative. Further details concerning the algorithms and criteria used for homology searches are provided below and are described in the pertinent literature highlighted by the citations provided below.

It will be appreciated that other criteria can be used to generate more inclusive and more 10 exclusive listings of the types set out in the tables. As those of skill will appreciate, narrow and broad searches both are useful. Thus, a skilled artisan can readily identify ORFs in contigs of the *T. pallidum* genome other than those specified for Tables 1-3, such as ORFs which are overlapping or encoded by the opposite strand of an identified ORF in addition to those ascertainable using the computer-based systems of the present invention.

15 As used herein, an "expression modulating fragment," EMF, means a series of nucleotide molecules which modulates the expression of an operably linked ORF or EMF.

As used herein, a sequence is said to "modulate the expression of an operably linked sequence" when the expression of the sequence is altered by the presence of the EMF. EMFs include, but are not limited to, promoters, and promoter modulating sequences (inducible 20 elements). One class of EMFs are fragments which induce the expression of an operably linked ORF in response to a specific regulatory factor or physiological event.

EMF sequences can be identified within the contigs of the *T. pallidum* genome by their proximity to the ORF IDs provided in Tables 1-3 and ORFs within each ORF ID. An intergenic segment, or a fragment of the intergenic segment, from about 10 to 200 nucleotides in length, 25 taken from any one of the ORFs of Tables 1-3 will modulate the expression of an operably linked ORF in a fashion similar to that found with the naturally linked ORF sequence. As used herein, an "intergenic segment" refers to fragments of the *T. pallidum* genome which are between two ORF(s) herein described. EMFs also can be identified using known EMFs as a target sequence or target motif in the computer-based systems of the present invention. Further, the two methods 30 can be combined and used together.

The presence and activity of an EMF can be confirmed using an EMF trap vector. An EMF trap vector contains a cloning site linked to a marker sequence. A marker sequence encodes an identifiable phenotype, such as antibiotic resistance or a complementing nutrition auxotrophic factor, which can be identified or assayed when the EMF trap vector is placed within an 35 appropriate host under appropriate conditions. As described above, a EMF will modulate the expression of an operably linked marker sequence. A more detailed discussion of various marker sequences is provided below. A sequence which is suspected as being an EMF is cloned in all three reading frames in one or more restriction sites upstream from the marker sequence in the EMF trap vector. The vector is then transformed into an appropriate host using known

procedures and the phenotype of the transformed host is examined under appropriate conditions. As described above, an EMF will modulate the expression of an operably linked marker sequence.

As used herein, a "diagnostic fragment," DF, means a series of nucleotide molecules which selectively hybridize to *T. pallidum* sequences. DFs can be readily identified by identifying unique sequences within contigs of the *T. pallidum* genome, such as by using well-known computer analysis software, and by generating and testing probes or amplification primers consisting of the DF sequence in an appropriate diagnostic format which determines amplification or hybridization selectivity.

The sequences falling within the scope of the present invention are not limited to the specific sequences herein described, but also include allelic and species variations thereof. Allelic and species variations can be routinely determined by comparing the polynucleotide sequences provided in SEQ ID NOS:1-744, ORF IDs and ORFs within, a representative fragment thereof, or a nucleotide sequence at least 99% and preferably 99.9% identical to said polynucleotide sequences, with a sequence from another isolate of the same species. Furthermore, to accommodate codon variability, the invention includes nucleic acid molecules coding for the same amino acid sequences as do the specific ORFs disclosed herein. In other words, in the coding region of an ORF, substitution of one codon for another which encodes the same amino acid is expressly contemplated.

Any specific sequence disclosed herein can be readily screened for errors by resequencing a particular fragment, such as an ORF, in both directions (*i.e.*, sequence both strands). Alternatively, error screening can be performed by sequencing corresponding polynucleotides of *T. pallidum* origin isolated by using part or all of the fragments in question as a probe or primer.

Each of the ORFs of the *T. pallidum* genome within the ORF IDs of Tables 1, 2 and 3, and the EMFs found 5' to the ORFs, can be used as polynucleotide reagents in numerous ways. For example, the sequences can be used as diagnostic probes or diagnostic amplification primers to detect the presence of a specific microbe in a sample, particularly *T. pallidum*. Especially preferred in this regard are ORFs such as those of Table 3, which do not match previously characterized sequences from other organisms and thus are most likely to be highly selective for *T. pallidum*. Also particularly preferred are ORFs that can be used to distinguish between strains of *T. pallidum*, particularly those that distinguish medically important strain, such as drug-resistant strains.

In addition, the fragments of the present invention, as broadly described, can be used to control gene expression through triple helix formation or antisense DNA or RNA, both of which methods are based on the binding of a polynucleotide sequence to DNA or RNA. Triple helix-formation optimally results in a shut-off of RNA transcription from DNA, while antisense RNA hybridization blocks translation of an mRNA molecule into polypeptide. Information from the sequences of the present invention can be used to design antisense and triple helix-forming oligonucleotides. Polynucleotides suitable for use in these methods are usually 20 to 40 bases in

length and are designed to be complementary to a region of the gene involved in transcription, for triple-helix formation, or to the mRNA itself, for antisense inhibition. Both techniques have been demonstrated to be effective in model systems, and the requisite techniques are well known and involve routine procedures. Triple helix techniques are discussed in, for example, Lee *et al.*,

- 5 *Nucl. Acids Res.* 6:3073 (1979); Cooney *et al.*, *Science* 241:456 (1988); and Dervan *et al.*, *Science* 251:1360 (1991). Antisense techniques in general are discussed in, for instance, Okano, *J. Neurochem.* 56:560 (1991) and *Oligodeoxynucleotides as Antisense Inhibitors of Gene Expression*, CRC Press, Boca Raton, FL (1988).

The present invention further provides recombinant constructs comprising one or more 10 fragments of the *T. pallidum* genomic fragments and contigs of the present invention. Certain preferred recombinant constructs of the present invention comprise a vector, such as a plasmid or viral vector, into which a fragment of the *T. pallidum* genome has been inserted, in a forward or reverse orientation. In the case of a vector comprising one of the ORFs of the present invention, the vector may further comprise regulatory sequences, including for example, a promoter, 15 operably linked to the ORF. For vectors comprising the EMFs of the present invention, the vector may further comprise a marker sequence or heterologous ORF operably linked to the EMF.

Large numbers of suitable vectors and promoters are known to those of skill in the art and are commercially available for generating the recombinant constructs of the present invention.

- 20 The following vectors are provided by way of example. Useful bacterial vectors include phagescript, PsiX174, pBluescript SK, pBS KS, pNH8a, pNH16a, pNH18a, pNH46a (available from Stratagene); pTrc99A, pKK223-3, pKK233-3, pDR540, pRIT5 (available from Pharmacia). Useful eukaryotic vectors include pWLneo, pSV2cat, pOG44, pXT1, pSG (available from Stratagene) pSVK3, pBPV, pMSG, pSVL (available from Pharmacia).

25 Promoter regions can be selected from any desired gene using CAT (chloramphenicol transferase) vectors or other vectors with selectable markers. Two appropriate vectors are pKK232-8 and pCM7. Particular named bacterial promoters include lacI, lacZ, T3, T7, gpt, lambda PR, and trc. Eukaryotic promoters include CMV immediate early, HSV thymidine kinase, early and late SV40, LTRs from retrovirus, and mouse metallothionein-I. Selection of 30 the appropriate vector and promoter is well within the level of ordinary skill in the art.

The present invention further provides host cells containing any one of the isolated 35 fragments of the *T. pallidum* genomic fragments and contigs of the present invention, wherein the fragment has been introduced into the host cell using known methods. The host cell can be a higher eukaryotic host cell, such as a mammalian cell, a lower eukaryotic host cell, such as a yeast cell, or a prokaryotic cell, such as a bacterial cell.

A polynucleotide of the present invention, such as a recombinant construct comprising an ORF of the present invention, may be introduced into the host by a variety of well established techniques that are standard in the art, such as calcium phosphate transfection, DEAE, dextran

mediated transfection and electroporation, which are described in, for instance, Davis, L. *et al.*, **BASIC METHODS IN MOLECULAR BIOLOGY** (1986).

A host cell containing one of the fragments of the *T. pallidum* genomic fragments and contigs of the present invention, can be used in conventional manners to produce the gene

- 5 product encoded by the isolated fragment (in the case of an ORF) or can be used to produce a heterologous protein under the control of the EMF.

The present invention further provides isolated polypeptides encoded by the nucleic acid fragments of the present invention or by degenerate variants of the nucleic acid fragments of the present invention. By "degenerate variant" is intended nucleotide fragments which differ from a 10 nucleic acid fragment of the present invention (*e.g.*, an ORF) by nucleotide sequence but, due to the degeneracy of the Genetic Code, encode an identical polypeptide sequence.

Preferred nucleic acid fragments of the present invention are the ORF IDs depicted in Tables 2 and 3 and the ORFs within which encode proteins.

A variety of methodologies known in the art can be utilized to obtain any one of the 15 isolated polypeptides or proteins of the present invention. At the simplest level, the amino acid sequence can be synthesized using commercially available peptide synthesizers. This is particularly useful in producing small peptides and fragments of larger polypeptides. Such short fragments as may be obtained most readily by synthesis are useful, for example, in generating antibodies against the native polypeptide, as discussed further below.

20 In an alternative method, the polypeptide or protein is purified from bacterial cells which naturally produce the polypeptide or protein. One skilled in the art can readily employ well-known methods for isolating polypeptides and proteins to isolate and purify polypeptides or proteins of the present invention produced naturally by a bacterial strain, or by other methods. Methods for isolation and purification that can be employed in this regard include, but are not 25 limited to, immunochromatography, HPLC, size-exclusion chromatography, ion-exchange chromatography, and immuno-affinity chromatography.

The polypeptides and proteins of the present invention also can be purified from cells which have been altered to express the desired polypeptide or protein. As used herein, a cell is said to be altered to express a desired polypeptide or protein when the cell, through genetic 30 manipulation, is made to produce a polypeptide or protein which it normally does not produce or which the cell normally produces at a lower level. Those skilled in the art can readily adapt procedures for introducing and expressing either recombinant or synthetic sequences into eukaryotic or prokaryotic cells in order to generate a cell which produces one of the polypeptides or proteins of the present invention.

35 The polypeptides of the present invention are preferably provided in an isolated form, and preferably are substantially purified. A recombinantly produced version of the *T. pallidum* polypeptide can be substantially purified by the one-step method described by Smith *et al.* (1988) Gene 67:31-40. Polypeptides of the invention also can be purified from natural or recombinant sources using antibodies directed against the polypeptides of the invention in methods which are

well known in the art of protein purification.

The invention further provides for isolated *T. pallidum* polypeptides comprising an amino acid sequence selected from the group including: (a) the amino acid sequence of a full-length *T. pallidum* polypeptide having the complete amino acid sequence from the first methionine codon to the termination codon of each sequence listed in SEQ ID NOS:1-744, wherein said termination codon is at the end of each SEQ ID NO: and said first methionine is the first methionine in frame with said termination codon; and (b) the amino acid sequence of a full-length *T. pallidum* polypeptide having the complete amino acid sequence in (a) excepting the N-terminal methionine.

The polypeptides of the present invention also include polypeptides having an amino acid sequence at least 80% identical, more preferably at least 90% identical, and still more preferably 95%, 96%, 97%, 98% or 99% identical to those described in (a) and (b) above.

The present invention is further directed to polynucleotides encoding portions or fragments of the amino acid sequences described herein as well as to portions or fragments of the isolated amino acid sequences described herein. Fragments include portions of the amino acid sequences described herein at least 5 contiguous amino acid in length and selected from any two integers, one of which representing an N-terminal position and another representing a C-terminal position. The initiation codon of the ORFs of the present invention is position 1. The initiation codon (positon 1) for purposes of the present invention is the first methionine codon of each ORF ID which is in frame with the termination codon at the end of each said sequence. Every combination of a N-terminal and C-terminal position that a fragment at least 5 contiguous amino acid residues in length could occupy, on any given ORF is included in the invention, i.e., from initiation codon up to the termination codon. "At least" means a fragment may be 5 contiguous amino acid residues in length or any integer between 5 and the number of residues in an ORF, minus 1. Therefore, included in the invention are contiguous fragments specified by any N-terminal and C-terminal positions of amino acid sequence set forth in SEQ ID NOS:1-744 or Tables 1-3 wherein the contiguous fragment is any integer between 5 and the number of residues in an ORF minus 1.

Further, the invention includes polypeptides comprising fragments specified by size, in amino acid residues, rather than by N-terminal and C-terminal positions. The invention includes any fragment size, in contiguous amino acid residues, selected from integers between 5 and the number of residues in an ORF, minus 1. Preferred sizes of contiguous polypeptide fragments include about 5 amino acid residues, about 10 amino acid residues, about 20 amino acid residues, about 30 amino acid residues, about 40 amino acid residues, about 50 amino acid residues, about 100 amino acid residues, about 200 amino acid residues, about 300 amino acid residues, and about 400 amino acid residues. The preferred sizes are, of course, meant to exemplify, not limit, the present invention as all size fragments representing any integer between 5 and the number of residues in a full length sequence minus 1 are included in the invention. The present invention also provides for the exclusion of any fragments specified by N-terminal and C-terminal positions or by size in amino acid residues as described above. Any number of fragments

specified by N-terminal and C-terminal positions or by size in amino acid residues as described above may be excluded.

The above fragments need not be active since they would be useful, for example, in immunoassays, in epitope mapping, epitope tagging, to generate antibodies to a particular portion of the protein, as vaccines, and as molecular weight markers.

Further polypeptides of the present invention include polypeptides which have at least 90% similarity, more preferably at least 95% similarity, and still more preferably at least 96%, 97%, 98% or 99% similarity to those described above.

A further embodiment of the invention relates to a polypeptide which comprises the amino acid sequence of a *T. pallidum* polypeptide having an amino acid sequence which contains at least one conservative amino acid substitution, but not more than 50 conservative amino acid substitutions, not more than 40 conservative amino acid substitutions, not more than 30 conservative amino acid substitutions, and not more than 20 conservative amino acid substitutions. Also provided are polypeptides which comprise the amino acid sequence of a *T. pallidum* polypeptide, having at least one, but not more than 10, 9, 8, 7, 6, 5, 4, 3, 2 or 1 conservative amino acid substitutions.

By a polypeptide having an amino acid sequence at least, for example, 95% "identical" to a query amino acid sequence of the present invention, it is intended that the amino acid sequence of the subject polypeptide is identical to the query sequence except that the subject polypeptide sequence may include up to five amino acid alterations per each 100 amino acids of the query amino acid sequence. In other words, to obtain a polypeptide having an amino acid sequence at least 95% identical to a query amino acid sequence, up to 5% of the amino acid residues in the subject sequence may be inserted, deleted, (indels) or substituted with another amino acid. These alterations of the reference sequence may occur at the amino or carboxy terminal positions of the reference amino acid sequence or anywhere between those terminal positions, interspersed either individually among residues in the reference sequence.

As a practical matter, whether any particular polypeptide is at least 90%, 95%, 96%, 97%, 98% or 99% identical to the ORF amino acid sequences encoded by the sequences of SEQ ID NOS:1-744, as described hererin, can be determined conventionally using known computer programs. A preferred method for determining the best overall match between a query sequence (a sequence of the present invention) and a subject sequence, also referred to as a global sequence alignment, can be determined using the FASTDB computer program based on the algorithm of Brutlag et al., (1990) Comp. App. Biosci. 6:237-245. In a sequence alignment the query and subject sequences are both amino acid sequences. The result of said global sequence alignment is in percent identity. Preferred parameters used in a FASTDB amino acid alignment are: Matrix=PAM 0, k-tuple=2, Mismatch Penalty=1, Joining Penalty=20, Randomization Group Length=0, Cutoff Score=1, Window Size=sequence length, Gap Penalty=5, Gap Size Penalty=0.05, Window Size=500 or the length of the subject amino acid sequence, whichever is shorter.

If the subject sequence is shorter than the query sequence due to N- or C-terminal deletions, not because of internal deletions, the results, in percent identity, must be manually corrected. This is because the FASTDB program does not account for N- and C-terminal truncations of the subject sequence when calculating global percent identity. For subject sequences truncated at the N- and C-termini, relative to the query sequence, the percent identity is corrected by calculating the number of residues of the query sequence that are N- and C-terminal of the subject sequence, which are not matched/aligned with a corresponding subject residue, as a percent of the total bases of the query sequence. Whether a residue is matched/aligned is determined by results of the FASTDB sequence alignment. This percentage is then subtracted from the percent identity, calculated by the above FASTDB program using the specified parameters, to arrive at a final percent identity score. This final percent identity score is what is used for the purposes of the present invention. Only residues to the N- and C-termini of the subject sequence, which are not matched/aligned with the query sequence, are considered for the purposes of manually adjusting the percent identity score. That is, only query amino acid residues outside the farthest N- and C-terminal residues of the subject sequence.

For example, a 90 amino acid residue subject sequence is aligned with a 100 residue query sequence to determine percent identity. The deletion occurs at the N-terminus of the subject sequence and therefore, the FASTDB alignment does not match/align with the first 10 residues at the N-terminus. The 10 unpaired residues represent 10% of the sequence (number of residues at the N- and C- termini not matched/total number of residues in the query sequence) so 10% is subtracted from the percent identity score calculated by the FASTDB program. If the remaining 90 residues were perfectly matched the final percent identity would be 90%. In another example, a 90 residue subject sequence is compared with a 100 residue query sequence. This time the deletions are internal so there are no residues at the N- or C-termini of the subject sequence which are not matched/aligned with the query. In this case the percent identity calculated by FASTDB is not manually corrected. Once again, only residue positions outside the N- and C-terminal ends of the subject sequence, as displayed in the FASTDB alignment, which are not matched/aligned with the query sequence are manually corrected. No other manual corrections are made for the purposes of the present invention.

The above polypeptide sequences are included irrespective of whether they have their normal biological activity. This is because even where a particular polypeptide molecule does not have biological activity, one of skill in the art would still know how to use the polypeptide, for instance, as a vaccine or to generate antibodies. Other uses of the polypeptides of the present invention that do not have *T. pallidum* activity include, *inter alia*, as epitope tags, in epitope mapping, and as molecular weight markers on SDS-PAGE gels or on molecular sieve gel filtration columns using methods known to those of skill in the art. As described below, the polypeptides of the present invention can also be used to raise polyclonal and monoclonal antibodies, which are useful in assays for detecting *T. pallidum* protein expression or as agonists and antagonists capable of enhancing or inhibiting *T. pallidum* protein

function. Further, such polypeptides can be used in the yeast two-hybrid system to "capture" *T. pallidum* protein binding proteins which are also candidate agonists and antagonists according to the present invention. See, e.g., Fields et al. (1989) Nature 340:245-246.

Any host/vector system can be used to express one or more of the ORFs of the present invention. These include, but are not limited to, eukaryotic hosts such as HeLa cells, CV-1 cell, COS cells, and Sf9 cells, as well as prokaryotic host such as *E. coli* and *B. subtilis*. The most preferred cells are those which do not normally express the particular polypeptide or protein or which expresses the polypeptide or protein at low natural level.

"Recombinant," as used herein, means that a polypeptide or protein is derived from recombinant (e.g., microbial or mammalian) expression systems. "Microbial" refers to recombinant polypeptides or proteins made in bacterial or fungal (e.g., yeast) expression systems. As a product, "recombinant microbial" defines a polypeptide or protein essentially free of native endogenous substances and unaccompanied by associated native glycosylation. Polypeptides or proteins expressed in most bacterial cultures, e.g., *E. coli*, will be free of glycosylation modifications; polypeptides or proteins expressed in yeast will have a glycosylation pattern different from that expressed in mammalian cells.

"Nucleotide sequence" refers to a heteropolymer of deoxyribonucleotides. Generally, DNA segments encoding the polypeptides and proteins provided by this invention are assembled from fragments of the *T. pallidum* genome and short oligonucleotide linkers, or from a series of oligonucleotides, to provide a synthetic gene which is capable of being expressed in a recombinant transcriptional unit comprising regulatory elements derived from a microbial or viral operon.

"Recombinant expression vehicle or vector" refers to a plasmid or phage or virus or vector, for expressing a polypeptide from a DNA (RNA) sequence. The expression vehicle can comprise a transcriptional unit comprising an assembly of (1) a genetic regulatory elements necessary for gene expression in the host, including elements required to initiate and maintain transcription at a level sufficient for suitable expression of the desired polypeptide, including, for example, promoters and, where necessary, an enhancer and a polyadenylation signal; (2) a structural or coding sequence which is transcribed into mRNA and translated into protein, and (3) appropriate signals to initiate translation at the beginning of the desired coding region and terminate translation at its end. Structural units intended for use in yeast or eukaryotic expression systems preferably include a leader sequence enabling extracellular secretion of translated protein by a host cell. Alternatively, where recombinant protein is expressed without a leader or transport sequence, it may include an N-terminal methionine residue. This residue may or may not be subsequently cleaved from the expressed recombinant protein to provide a final product.

"Recombinant expression system" means host cells which have stably integrated a recombinant transcriptional unit into chromosomal DNA or carry the recombinant transcriptional unit extra chromosomally. The cells can be prokaryotic or eukaryotic. Recombinant expression

systems as defined herein will express heterologous polypeptides or proteins upon induction of the regulatory elements linked to the DNA segment or synthetic gene to be expressed.

Mature proteins can be expressed in mammalian cells, yeast, bacteria, or other cells under the control of appropriate promoters. Cell-free translation systems can also be employed to

5 produce such proteins using RNAs derived from the DNA constructs of the present invention.

Appropriate cloning and expression vectors for use with prokaryotic and eukaryotic hosts are described in Sambrook *et al.*, *Molecular Cloning: A Laboratory Manual*, 2nd Edition, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York (1989), the disclosure of which is hereby incorporated by reference in its entirety.

10 Generally, recombinant expression vectors will include origins of replication and selectable markers permitting transformation of the host cell, *e.g.*, the ampicillin resistance gene of *E. coli* and *S. cerevisiae* TRP1 gene, and a promoter derived from a highly expressed gene to direct transcription of a downstream structural sequence. Such promoters can be derived from operons encoding glycolytic enzymes such as 3-phosphoglycerate kinase (PGK), alpha-factor,

15 acid phosphatase, or heat shock proteins, among others. The heterologous structural sequence is assembled in appropriate phase with translation initiation and termination sequences, and preferably, a leader sequence capable of directing secretion of translated protein into the periplasmic space or extracellular medium. Optionally, the heterologous sequence can encode a fusion protein including an N-terminal identification peptide imparting desired characteristics,

20 *e.g.*, stabilization or simplified purification of expressed recombinant product.

Useful expression vectors for bacterial use are constructed by inserting a structural DNA sequence encoding a desired protein together with suitable translation initiation and termination signals in operable reading phase with a functional promoter. The vector will comprise one or more phenotypic selectable markers and an origin of replication to ensure maintenance of the 25 vector and, when desirable, provide amplification within the host.

Suitable prokaryotic hosts for transformation include strains of *E. coli*, *B. subtilis*, *Salmonella typhimurium* and various species within the genera *Pseudomonas* and *Streptomyces*. Others may, also be employed as a matter of choice.

As a representative but non-limiting example, useful expression vectors for bacterial use 30 can comprise a selectable marker and bacterial origin of replication derived from commercially available plasmids comprising genetic elements of the well known cloning vector pBR322 (ATCC 37017). Such commercial vectors include, for example, pKK223-3 (available from Pharmacia Fine Chemicals, Uppsala, Sweden) and GEM 1 (available from Promega Biotec, Madison, WI, USA). These pBR322 "backbone" sections are combined with an appropriate 35 promoter and the structural sequence to be expressed.

Following transformation of a suitable host strain and growth of the host strain to an appropriate cell density, the selected promoter, where it is inducible, is derepressed or induced by appropriate means (*e.g.*, temperature shift or chemical induction) and cells are cultured for an additional period to provide for expression of the induced gene product. Thereafter cells are

typically harvested, generally by centrifugation, disrupted to release expressed protein, generally by physical or chemical means, and the resulting crude extract is retained for further purification.

Various mammalian cell culture systems can also be employed to express recombinant protein. Examples of mammalian expression systems include the COS-7 lines of monkey kidney fibroblasts, described in Gluzman, *Cell* 23:175 (1981), and other cell lines capable of expressing a compatible vector, for example, the C127, 3T3, CHO, HeLa and BHK cell lines.

Mammalian expression vectors will comprise an origin of replication, a suitable promoter and enhancer, and also any necessary ribosome binding sites, polyadenylation site, splice donor and acceptor sites, transcriptional termination sequences, and 5' flanking nontranscribed sequences. DNA sequences derived from the SV40 viral genome, for example, SV40 origin, early promoter, enhancer, splice, and polyadenylation sites may be used to provide the required nontranscribed genetic elements.

Recombinant polypeptides and proteins produced in bacterial culture is usually isolated by initial extraction from cell pellets, followed by one or more salting-out, aqueous ion exchange or size exclusion chromatography steps. Microbial cells employed in expression of proteins can be disrupted by any convenient method, including freeze-thaw cycling, sonication, mechanical disruption, or use of cell lysing agents. Protein refolding steps can be used, as necessary, in completing configuration of the mature protein. Finally, high performance liquid chromatography (HPLC) can be employed for final purification steps.

The present invention further includes isolated polypeptides, proteins and nucleic acid molecules which are substantially equivalent to those herein described. As used herein, substantially equivalent can refer both to nucleic acid and amino acid sequences, for example a mutant sequence, that varies from a reference sequence by one or more substitutions, deletions, or additions, the net effect of which does not result in an adverse functional dissimilarity between reference and subject sequences. For purposes of the present invention, sequences having equivalent biological activity, and equivalent expression characteristics are considered substantially equivalent. For purposes of determining equivalence, truncation of the mature sequence should be disregarded.

The invention further provides methods of obtaining homologs from other strains of *T. pallidum*, of the fragments of the *T. pallidum* genome of the present invention and homologs of the proteins encoded by the ORFs of the present invention. As used herein, a sequence or protein of *T. pallidum* is defined as a homolog of a fragment of the *T. pallidum* fragments or contigs or a protein encoded by one of the ORFs of the present invention, if it shares significant homology to one of the fragments of the *T. pallidum* genome of the present invention or a protein encoded by one of the ORFs of the present invention. Specifically, by using the sequence disclosed herein as a probe or as primers, and techniques such as PCR cloning and colony/plaque hybridization, one skilled in the art can obtain homologs.

As used herein, two nucleic acid molecules or proteins are said to "share significant homology" if the two contain regions which possess greater than 85% sequence (amino acid or

nucleic acid) homology. Preferred homologs in this regard are those with more than 90% homology. Especially preferred are those with 93% or more homology. Among especially preferred homologs those with 95% or more homology are particularly preferred. Very particularly preferred among these are those with 97% and even more particularly preferred

- 5 among those are homologs with 99% or more homology. The most preferred homologs among these are those with 99.9% homology or more. It will be understood that, among measures of homology, identity is particularly preferred in this regard.

Region specific primers or probes derived from the nucleotide sequence provided in SEQ ID NOS: 1-744 or from a nucleotide sequence at least 95%, particularly at least 99%, especially 10 at least 99.5% identical to a sequence of SEQ ID NOS: 1-744 can be used to prime DNA synthesis and PCR amplification, as well as to identify colonies containing cloned DNA encoding a homolog. Methods suitable to this aspect of the present invention are well known and have been described in great detail in many publications such as, for example, Innis *et al.*, *PCR Protocols*, Academic Press, San Diego, CA (1990)).

15 When using primers derived from SEQ ID NOS: 1-744 or from a nucleotide sequence having an aforementioned identity to a sequence of SEQ ID NOS: 1-744, one skilled in the art will recognize that by employing high stringency conditions (e.g., annealing at 50-60°C in 6X SSPC and 50% formamide, and washing at 50- 65°C in 0.5X SSPC) only sequences which are greater than 75% homologous to the primer will be amplified. By employing lower stringency 20 conditions (e.g., hybridizing at 35-37°C in 5X SSPC and 40-45% formamide, and washing at 42°C in 0.5X SSPC), sequences which are greater than 40-50% homologous to the primer will also be amplified.

When using DNA probes derived from SEQ ID NOS: 1-744, or from a nucleotide sequence having an aforementioned identity to a sequence of SEQ ID NOS: 1-744 , for 25 colony/plaque hybridization, one skilled in the art will recognize that by employing high stringency conditions (e.g., hybridizing at 50- 65°C in 5X SSPC and 50% formamide, and washing at 50- 65°C in 0.5X SSPC), sequences having regions which are greater than 90% homologous to the probe can be obtained, and that by employing lower stringency conditions (e.g., hybridizing at 35-37°C in 5X SSPC and 40-45% formamide, and washing at 42°C in 0.5X 30 SSPC), sequences having regions which are greater than 35-45% homologous to the probe will be obtained.

Any organism can be used as the source for homologs of the present invention so long as the organism naturally expresses such a protein or contains genes encoding the same. The most preferred organism for isolating homologs are bacteria which are closely related to *T. pallidum*.

35

ILLUSTRATIVE USES OF COMPOSITIONS OF THE INVENTION

Each ORF corresponding to the ORF IDs provided in Tables 1 and 2 is identified with a function by homology to a known gene or polypeptide. As a result, one skilled in the art can use

the polypeptides of the present invention for commercial, therapeutic and industrial purposes consistent with the type of putative identification of the polypeptide. Such identifications permit one skilled in the art to use the *T. pallidum* ORFs in a manner similar to the known type of sequences for which the identification is made; for example, to ferment a particular sugar source or to produce a particular metabolite. A variety of reviews illustrative of this aspect of the invention are available, including the following reviews on the industrial use of enzymes, for example, BIOCHEMICAL ENGINEERING AND BIOTECHNOLOGY HANDBOOK, 2nd Ed., MacMillan Publications, Ltd. NY (1991) and BIOCATALYSTS IN ORGANIC SYNTHESSES, Tramper *et al.*, Eds., Elsevier Science Publishers, Amsterdam, The Netherlands (1985). A variety of exemplary uses that illustrate this and similar aspects of the present invention are discussed below.

1. Biosynthetic Enzymes

Open reading frames encoding proteins involved in mediating the catalytic reactions involved in intermediary and macromolecular metabolism, the biosynthesis of small molecules, cellular processes and other functions includes enzymes involved in the degradation of the intermediary products of metabolism, enzymes involved in central intermediary metabolism, enzymes involved in respiration, both aerobic and anaerobic, enzymes involved in fermentation, enzymes involved in ATP proton motor force conversion, enzymes involved in broad regulatory function, enzymes involved in amino acid synthesis, enzymes involved in nucleotide synthesis, enzymes involved in cofactor and vitamin synthesis, can be used for industrial biosynthesis.

The various metabolic pathways present in *T. pallidum* can be identified based on absolute nutritional requirements as well as by examining the various enzymes identified in Table 1-3 and SEQ ID NOS:1-744.

Of particular interest are polypeptides involved in the degradation of intermediary metabolites as well as non-macromolecular metabolism. Such enzymes include amylases, glucose oxidases, and catalase.

Proteolytic enzymes are another class of commercially important enzymes. Proteolytic enzymes find use in a number of industrial processes including the processing of flax and other vegetable fibers, in the extraction, clarification and depectinization of fruit juices, in the extraction of vegetables' oil and in the maceration of fruits and vegetables to give unicellular fruits. A detailed review of the proteolytic enzymes used in the food industry is provided in Rombouts *et al.*, *Symbiosis* 21:79 (1986) and Voragen *et al.* in *Biocatalysts In Agricultural Biotechnology*, Whitaker *et al.*, Eds., *American Chemical Society Symposium Series* 389:93 (1989).

The metabolism of sugars is an important aspect of the primary metabolism of *T. pallidum*. Enzymes involved in the degradation of sugars, such as, particularly, glucose, galactose, fructose and xylose, can be used in industrial fermentation. Some of the important sugar transforming enzymes, from a commercial viewpoint, include sugar isomerases such as glucose isomerase. Other metabolic enzymes have found commercial use such as glucose

oxidases which produces ketogulonic acid (KGA). KGA is an intermediate in the commercial production of ascorbic acid using the Reichstein's procedure, as described in Krueger *et al.*, *Biotechnology 6(A)*, Rhine *et al.*, Eds., Verlag Press, Weinheim, Germany (1984).

Glucose oxidase (GOD) is commercially available and has been used in purified form as

- 5 well as in an immobilized form for the deoxygenation of beer. See, for instance, Hartmeir *et al.*, *Biotechnology Letters 1:21* (1979). The most important application of GOD is the industrial scale fermentation of gluconic acid. Market for gluconic acids which are used in the detergent, textile, leather, photographic, pharmaceutical, food, feed and concrete industry, as described, for example, in Bigelis *et al.*, beginning on page 357 in *GENE MANIPULATIONS AND FUNGI*;
- 10 Bennett *et al.*, Eds., Academic Press, New York (1985). In addition to industrial applications, GOD has found applications in medicine for quantitative determination of glucose in body fluids recently in biotechnology for analyzing syrups from starch and cellulose hydrolysates. This application is described in Owusu *et al.*, *Biochem. et Biophysica. Acta.* 872:83 (1986), for instance.

- 15 The main sweetener used in the world today is sugar which comes from sugar beets and sugar cane. In the field of industrial enzymes, the glucose isomerase process shows the largest expansion in the market today. Initially, soluble enzymes were used and later immobilized enzymes were developed (Krueger *et al.*, *Biotechnology, The Textbook of Industrial Microbiology*, Sinauer Associated Incorporated, Sunderland, Massachusetts (1990)). Today, the 20 use of glucose- produced high fructose syrups is by far the largest industrial business using immobilized enzymes. A review of the industrial use of these enzymes is provided by Jorgensen, *Starch 40:307* (1988).

- 25 Proteinases, such as alkaline serine proteinases, are used as detergent additives and thus represent one of the largest volumes of microbial enzymes used in the industrial sector. Because of their industrial importance, there is a large body of published and unpublished information regarding the use of these enzymes in industrial processes. (See Faultman *et al.*, *Acid Proteases Structure Function and Biology*, Tang, J., ed., Plenum Press, New York (1977) and Godfrey *et al.*, *Industrial Enzymes*, MacMillan Publishers, Surrey, UK (1983) and Hepner *et al.*, Report Industrial Enzymes by 1990, Hel Hepner & Associates, London (1986)).

- 30 Another class of commercially usable proteins of the present invention are the microbial lipases, described by, for instance, Macrae *et al.*, *Philosophical Transactions of the Chiral Society of London 310:227* (1985) and Poserke, *Journal of the American Oil Chemist Society 61:1758* (1984). A major use of lipases is in the fat and oil industry for the production of neutral glycerides using lipase catalyzed inter-esterification of readily available triglycerides. Application 35 of lipases include the use as a detergent additive to facilitate the removal of fats from fabrics in the course of the washing procedures.

The use of enzymes, and in particular microbial enzymes, as catalyst for key steps in the synthesis of complex organic molecules is gaining popularity at a great rate. One area of great interest is the preparation of chiral intermediates. Preparation of chiral intermediates is of interest

to a wide range of synthetic chemists particularly those scientists involved with the preparation of new pharmaceuticals, agrochemicals, fragrances and flavors. (See Davies *et al.*, *Recent Advances in the Generation of Chiral Intermediates Using Enzymes*, CRC Press, Boca Raton, Florida (1990)). The following reactions catalyzed by enzymes are of interest to organic

- 5 chemists: hydrolysis of carboxylic acid esters, phosphate esters, amides and nitriles, esterification reactions, trans-esterification reactions, synthesis of amides, reduction of alkanones and oxoalkanates, oxidation of alcohols to carbonyl compounds, oxidation of sulfides to sulfoxides, and carbon bond forming reactions such as the aldol reaction.

When considering the use of an enzyme encoded by one of the ORFs of the present
10 invention for biotransformation and organic synthesis it is sometimes necessary to consider the respective advantages and disadvantages of using a microorganism as opposed to an isolated enzyme. Pros and cons of using a whole cell system on the one hand or an isolated partially purified enzyme on the other hand, has been described in detail by Bud *et al.*, *Chemistry in Britain* (1987), p. 127.

- 15 Amino transferases, enzymes involved in the biosynthesis and metabolism of amino acids, are useful in the catalytic production of amino acids. The advantages of using microbial based enzyme systems is that the amino transferase enzymes catalyze the stereo-selective synthesis of only L-amino acids and generally possess uniformly high catalytic rates. A description of the use of amino transferases for amino acid production is provided by Roselle-
20 David, *Methods of Enzymology* 136:479 (1987).

Another category of useful proteins encoded by the ORFs of the present invention include enzymes involved in nucleic acid synthesis, repair, and recombination.

2. Generation of Antibodies

- 25 As described here, the proteins of the present invention, as well as homologs thereof, can be used in a variety of procedures and methods known in the art which are currently applied to other proteins. The proteins of the present invention can further be used to generate an antibody which selectively binds the protein.

- 30 *T. pallidum* protein-specific antibodies for use in the present invention can be raised against the intact *T. pallidum* protein or an antigenic polypeptide fragment thereof, which may be presented together with a carrier protein, such as an albumin, to an animal system (such as rabbit or mouse) or, if it is long enough (at least about 25 amino acids), without a carrier.

- 35 As used herein, the term "antibody" (Ab) or "monoclonal antibody" (Mab) is meant to include intact molecules, single chain whole antibodies, and antibody fragments. Antibody fragments of the present invention include Fab and F(ab')2 and other fragments including single-chain Fvs (scFv) and disulfide-linked Fvs (sdFv). Also included in the present invention are chimeric and humanized monoclonal antibodies and polyclonal antibodies specific for the polypeptides of the present invention. The antibodies of the present invention may be prepared by any of a variety of methods. For example, cells expressing a polypeptide of the present

invention or an antigenic fragment thereof can be administered to an animal in order to induce the production of sera containing polyclonal antibodies. For example, a preparation of *T. pallidum* polypeptide or fragment thereof is prepared and purified to render it substantially free of natural contaminants. Such a preparation is then introduced into an animal in order to produce 5 polyclonal antisera of greater specific activity.

In a preferred method, the antibodies of the present invention are monoclonal antibodies or binding fragments thereof. Such monoclonal antibodies can be prepared using hybridoma technology. See, e.g., Harlow et al., ANTIBODIES: A LABORATORY MANUAL, (Cold Spring Harbor Laboratory Press, 2nd ed. 1988); Hammerling, et al., in: MONOCLONAL 10 ANTIBODIES AND T-CELL HYBRIDOMAS 563-681 (Elsevier, N.Y., 1981). Fab and F(ab')2 fragments may be produced by proteolytic cleavage, using enzymes such as papain (to produce Fab fragments) or pepsin (to produce F(ab')2 fragments). Alternatively, *T. pallidum* polypeptide-binding fragments, chimeric, and humanized antibodies can be produced through the application of recombinant DNA technology or through synthetic chemistry using methods 15 known in the art.

Alternatively, additional antibodies capable of binding to the polypeptide antigen of the present invention may be produced in a two-step procedure through the use of anti-idiotypic antibodies. Such a method makes use of the fact that antibodies are themselves antigens, and that, therefore, it is possible to obtain an antibody which binds to a second antibody. In 20 accordance with this method, *T. pallidum* polypeptide-specific antibodies are used to immunize an animal, preferably a mouse. The splenocytes of such an animal are then used to produce hybridoma cells, and the hybridoma cells are screened to identify clones which produce an antibody whose ability to bind to the *T. pallidum* polypeptide-specific antibody can be blocked by the *T. pallidum* polypeptide antigen. Such antibodies comprise anti-idiotypic antibodies to 25 the *T. pallidum* polypeptide-specific antibody and can be used to immunize an animal to induce formation of further *T. pallidum* polypeptide-specific antibodies.

Antibodies and fragement thereof of the present invention may be described by the portion of a polypeptide of the present invention recognized or specifically bound by the antibody. Antibody binding fragement of a polypeptide of the present invention may be 30 described or specified in the same manner as for polypeptide framents discussed above., i.e., by N-terminal and C-terminal positions or by size in contiguous amino acid residues. Any number of antibody binding framents, of a polypeptide of the present invention, specified by N-terminal and C-terminal positions or by size in amino acid residues, as described above, may also be excluded from the present invention. Therefore, the present invention includes antibodies the 35 specifically bind a particullarly discribed frament of a polypeptide of the present invention and allows for the exclusion of the same.

Antibodies and framents thereof of the present invention may also be described or specified in terms of their cross-reactivity. Antibodies and framents that do not bind polypeptides of any other species of *Borrelia* other than *T. pallidum* are included in the present invention. Likewise,

antibodies and fragement that bind only species of *Borrelia*, i.e. antibodies and fragement that do not bind bacteria from any genus other than *Borrelia*, are included in the present invention.

The present invention further provides the above-described antibodies in detectably labelled form. Antibodies can be detectably labelled through the use of radioisotopes, affinity labels (such as biotin, avidin, etc.), enzymatic labels (such as horseradish peroxidase, alkaline phosphatase, etc.) fluorescent labels (such as FITC or rhodamine, etc.), paramagnetic atoms, etc. Procedures for accomplishing such labeling are well-known in the art, for example see Sternberger *et al.*, *J. Histochem. Cytochem.* 18:315 (1970); Bayer, E. A. *et al.*, *Meth. Enzym.* 62:308 (1979); Engval, E. *et al.*, *Immunol.* 109:129 (1972); Goding, J. W., *J. Immunol. Meth.* 13:215 (1976)).

The labeled antibodies of the present invention can be used for *in vitro*, *in vivo*, and *in situ* assays to identify cells or tissues in which a fragment of the *T. pallidum* genome is expressed.

The present invention further provides the above-described antibodies immobilized on a solid support. Examples of such solid supports include plastics such as polycarbonate, complex carbohydrates such as agarose and sepharose, acrylic resins and such as polyacrylamide and latex beads. Techniques for coupling antibodies to such solid supports are well known in the art (Weir, D. M. *et al.*, "Handbook of Experimental Immunology" 4th Ed., Blackwell Scientific Publications, Oxford, England, Chapter 10 (1986); Jacoby, W. D. *et al.*, *Meth. Enzym.* 34 Academic Press, N. Y. (1974)). The immobilized antibodies of the present invention can be used for *in vitro*, *in vivo*, and *in situ* assays as well as for immunoaffinity purification of the proteins of the present invention.

3. Epitope-Bearing Portions

In another aspect, the invention provides peptides and polypeptides comprising epitope-bearing portions of the *T. pallidum* polypeptides of the present invention. These epitopes are immunogenic or antigenic epitopes of the polypeptides of the present invention. An "immunogenic epitope" is defined as a part of a protein that elicits an antibody response when the whole protein or polypeptide is the immunogen. These immunogenic epitopes are believed to be confined to a few loci on the molecule. On the other hand, a region of a protein molecule to which an antibody can bind is defined as an "antigenic determinant" or "antigenic epitope." The number of immunogenic epitopes of a protein generally is less than the number of antigenic epitopes. See, e.g., Geysen, et al. (1983) *Proc. Natl. Acad. Sci. USA* 81:3998- 4002. Amino acid residues comprising antigenic epitopes may be determined by algorithms such as the Jameson-Wolf analysis or similar algorithms or by *in vivo* testing for an antigenic response using the methods described herein or those known in the art.

As to the selection of peptides or polypeptides bearing an antigenic epitope (i.e., that contain a region of a protein molecule to which an antibody can bind), it is well known in that art that relatively short synthetic peptides that mimic part of a protein sequence are routinely capable

of eliciting an antiserum that reacts with the partially mimicked protein. *See, e.g.*, Sutcliffe, et al., (1983) *Science* 219:660-666. Peptides capable of eliciting protein-reactive sera are frequently represented in the primary sequence of a protein, can be characterized by a set of simple chemical rules, and are confined neither to immunodominant regions of intact proteins (i.e., immunogenic epitopes) nor to the amino or carboxyl terminals. Peptides that are extremely hydrophobic and those of six or fewer residues generally are ineffective at inducing antibodies that bind to the mimicked protein; longer, peptides, especially those containing proline residues, usually are effective. *See, Sutcliffe, et al., supra*, p. 661. For instance, 18 of 20 peptides designed according to these guidelines, containing 8-39 residues covering 75% of the sequence of the influenza virus hemagglutinin HA1 polypeptide chain, induced antibodies that reacted with the HA1 protein or intact virus; and 12/12 peptides from the MuLV polymerase and 18/18 from the rabies glycoprotein induced antibodies that precipitated the respective proteins.

Antigenic epitope-bearing peptides and polypeptides of the invention are therefore useful to raise antibodies, including monoclonal antibodies, that bind specifically to a polypeptide of the invention. Thus, a high proportion of hybridomas obtained by fusion of spleen cells from donors immunized with an antigen epitope-bearing peptide generally secrete antibody reactive with the native protein. *See Sutcliffe, et al., supra*, p. 663. The antibodies raised by antigenic epitope-bearing peptides or polypeptides are useful to detect the mimicked protein, and antibodies to different peptides may be used for tracking the fate of various regions of a protein precursor which undergoes post-translational processing. The peptides and anti-peptide antibodies may be used in a variety of qualitative or quantitative assays for the mimicked protein, for instance in competition assays since it has been shown that even short peptides (*e.g.*, about 9 amino acids) can bind and displace the larger peptides in immunoprecipitation assays. *See, e.g.*, Wilson, et al., (1984) *Cell* 37:767-778. The anti-peptide antibodies of the invention also are useful for purification of the mimicked protein, for instance, by adsorption chromatography using methods known in the art.

Antigenic epitope-bearing peptides and polypeptides of the invention designed according to the above guidelines preferably contain a sequence of at least seven, more preferably at least nine and most preferably between about 10 to about 50 amino acids (i.e. any integer between 7 and 50) contained within the amino acid sequence of a polypeptide of the invention. However, peptides or polypeptides comprising a larger portion of an amino acid sequence of a polypeptide of the invention, containing about 50 to about 100 amino acids, or any length up to and including the entire amino acid sequence of a polypeptide of the invention, also are considered epitope-bearing peptides or polypeptides of the invention and also are useful for inducing antibodies that react with the mimicked protein. Preferably, the amino acid sequence of the epitope-bearing peptide is selected to provide substantial solubility in aqueous solvents (*i.e.*, the sequence includes relatively hydrophilic residues and highly hydrophobic sequences are preferably avoided); and sequences containing proline residues are particularly preferred.

The epitope-bearing peptides and polypeptides of the present invention may be produced by any conventional means for making peptides or polypeptides including recombinant means using nucleic acid molecules of the invention. For instance, an epitope-bearing amino acid sequence of the present invention may be fused to a larger polypeptide which acts as a carrier during recombinant production and purification, as well as during immunization to produce anti-peptide antibodies. Epitope-bearing peptides also may be synthesized using known methods of chemical synthesis. For instance, Houghten has described a simple method for synthesis of large numbers of peptides, such as 10-20 mg of 248 different 13 residue peptides representing single amino acid variants of a segment of the HA1 polypeptide which were prepared and characterized (by ELISA-type binding studies) in less than four weeks (Houghten, R. A. Proc. Natl. Acad. Sci. USA 82:5131-5135 (1985)). This "Simultaneous Multiple Peptide Synthesis (SMPS)" process is further described in U.S. Patent No. 4,631,211 to Houghten and coworkers (1986). In this procedure the individual resins for the solid-phase synthesis of various peptides are contained in separate solvent-permeable packets, enabling the optimal use of the many identical repetitive steps involved in solid-phase methods. A completely manual procedure allows 500-1000 or more syntheses to be conducted simultaneously (Houghten et al. (1985) Proc. Natl. Acad. Sci. 82:5131-5135 at 5134).

Epitope-bearing peptides and polypeptides of the invention are used to induce antibodies according to methods well known in the art. See, e.g., Sutcliffe, et al., *supra*; Wilson, et al., *supra*; and Bittle, et al. (1985) J. Gen. Virol. 66:2347-2354. Generally, animals may be immunized with free peptide; however, anti-peptide antibody titer may be boosted by coupling of the peptide to a macromolecular carrier, such as keyhole limpet hemacyanin (KLH) or tetanus toxoid. For instance, peptides containing cysteine may be coupled to carrier using a linker such as m-maleimidobenzoyl-N-hydroxysuccinimide ester (MBS), while other peptides may be coupled to carrier using a more general linking agent such as glutaraldehyde. Animals such as rabbits, rats and mice are immunized with either free or carrier-coupled peptides, for instance, by intraperitoneal and/or intradermal injection of emulsions containing about 100 µg peptide or carrier protein and Freund's adjuvant. Several booster injections may be needed, for instance, at intervals of about two weeks, to provide a useful titer of anti-peptide antibody which can be detected, for example, by ELISA assay using free peptide adsorbed to a solid surface. The titer of anti-peptide antibodies in serum from an immunized animal may be increased by selection of anti-peptide antibodies, for instance, by adsorption to the peptide on a solid support and elution of the selected antibodies according to methods well known in the art.

Immunogenic epitope-bearing peptides of the invention, i.e., those parts of a protein that elicit an antibody response when the whole protein is the immunogen, are identified according to methods known in the art. For instance, Geysen, et al., *supra*, discloses a procedure for rapid concurrent synthesis on solid supports of hundreds of peptides of sufficient purity to react in an ELISA. Interaction of synthesized peptides with antibodies is then easily detected without removing them from the support. In this manner a peptide bearing an immunogenic epitope of a

desired protein may be identified routinely by one of ordinary skill in the art. For instance, the immunologically important epitope in the coat protein of foot-and-mouth disease virus was located by Geysen *et al. supra* with a resolution of seven amino acids by synthesis of an overlapping set of all 208 possible hexapeptides covering the entire 213 amino acid sequence of the protein. Then, a complete replacement set of peptides in which all 20 amino acids were substituted in turn at every position within the epitope were synthesized, and the particular amino acids conferring specificity for the reaction with antibody were determined. Thus, peptide analogs of the epitope-bearing peptides of the invention can be made routinely by this method. U.S. Patent No. 4,708,781 to Geysen (1987) further describes this method of identifying a peptide bearing an immunogenic epitope of a desired protein.

Further still, U.S. Patent No. 5,194,392, to Geysen (1990), describes a general method of detecting or determining the sequence of monomers (amino acids or other compounds) which is a topological equivalent of the epitope (*i.e.*, a "mimotope") which is complementary to a particular paratope (antigen binding site) of an antibody of interest. More generally, U.S. Patent No. 4,433,092, also to Geysen (1989), describes a method of detecting or determining a sequence of monomers which is a topographical equivalent of a ligand which is complementary to the ligand binding site of a particular receptor of interest. Similarly, U.S. Patent No. 5,480,971 to Houghten, R. A. *et al.* (1996) discloses linear C₁-C₇-alkyl peralkylated oligopeptides and sets and libraries of such peptides, as well as methods for using such oligopeptide sets and libraries for determining the sequence of a peralkylated oligopeptide that preferentially binds to an acceptor molecule of interest. Thus, non-peptide analogs of the epitope-bearing peptides of the invention also can be made routinely by these methods. The entire disclosure of each document cited in this section on "Polypeptides and Fragments" is hereby incorporated herein by reference.

As one of skill in the art will appreciate, the polypeptides of the present invention and the epitope-bearing fragments thereof described above can be combined with parts of the constant domain of immunoglobulins (IgG), resulting in chimeric polypeptides. These fusion proteins facilitate purification and show an increased half-life *in vivo*. This has been shown, *e.g.*, for chimeric proteins consisting of the first two domains of the human CD4-polypeptide and various domains of the constant regions of the heavy or light chains of mammalian immunoglobulins. (EPA 0,394,827; Traunecker *et al.* (1988) Nature 331:84-86. Fusion proteins that have a disulfide-linked dimeric structure due to the IgG part can also be more efficient in binding and neutralizing other molecules than a monomeric *T. pallidum* polypeptide or fragment thereof alone. See Fountoulakis *et al.* (1995) J. Biochem. 270:3958-3964. Nucleic acids encoding the above epitopes of *T. pallidum* polypeptides can also be recombined with a gene of interest as an epitope tag to aid in detection and purification of the expressed polypeptide.

3. Diagnostic Assays and Kits

The present invention further relates to methods for assaying *Borrelia* infection in an animal by detecting the expression of genes encoding *Borrelia* polypeptides of the present invention. The methods comprise analyzing tissue or body fluid from the animal for

- 5 *Borrelia*-specific antibodies, nucleic acids, or proteins. Analysis of nucleic acid specific to *Borrelia* is assayed by PCR or hybridization techniques using nucleic acid sequences of the present invention as either hybridization probes or primers. See, e.g., Sambrook et al. Molecular cloning: A Laboratory Manual (Cold Spring Harbor Laboratory Press, 2nd ed., 1989, page 54 reference); Eremeeva et al. (1994) J. Clin. Microbiol. 32:803-810 (describing
10 differentiation among spotted fever group *Rickettsiae* species by analysis of restriction fragment length polymorphism of PCR-amplified DNA) and Chen et al. 1994 J. Clin. Microbiol. 32:589-595 (detecting *T. pallidum* nucleic acids via PCR).

Where diagnosis of a disease state related to infection with *Borrelia* has already been made, the present invention is useful for monitoring progression or regression of the disease state
15 whereby patients exhibiting enhanced *Borrelia* gene expression will experience a worse clinical outcome relative to patients expressing these gene(s) at a lower level.

By "biological sample" is intended any biological sample obtained from an animal, cell line, tissue culture, or other source which contains *Borrelia* polypeptide, mRNA, or DNA. Biological samples include body fluids (such as saliva, blood, plasma, urine, mucus, synovial
20 fluid, etc.) tissues (such as muscle, skin, and cartilage) and any other biological source suspected of containing *Borrelia* polypeptides or nucleic acids. Methods for obtaining biological samples such as tissue are well known in the art.

The present invention is useful for detecting diseases related to *Borrelia* infections in animals. Preferred animals include monkeys, apes, cats, dogs, birds, cows, pigs, mice, horses,
25 rabbits and humans. Particularly preferred are humans.

Total RNA can be isolated from a biological sample using any suitable technique such as the single-step guanidinium-thiocyanate-phenol-chloroform method described in Chomczynski et al. (1987) Anal. Biochem. 162:156-159. mRNA encoding *Borrelia* polypeptides having sufficient homology to the nucleic acid sequences identified in SEQ ID NOS:1-744 to allow for
30 hybridization between complementary sequences are then assayed using any appropriate method. These include Northern blot analysis, S1 nuclease mapping, the polymerase chain reaction (PCR), reverse transcription in combination with the polymerase chain reaction (RT-PCR), and reverse transcription in combination with the ligase chain reaction (RT-LCR).

Northern blot analysis can be performed as described in Harada et al. (1990) Cell
35 63:303-312. Briefly, total RNA is prepared from a biological sample as described above. For the Northern blot, the RNA is denatured in an appropriate buffer (such as glyoxal/dimethyl sulfoxide/sodium phosphate buffer), subjected to agarose gel electrophoresis, and transferred onto a nitrocellulose filter. After the RNAs have been linked to the filter by a UV linker, the filter is prehybridized in a solution containing formamide, SSC, Denhardt's solution, denatured

salmon sperm, SDS, and sodium phosphate buffer. A *T. pallidum* polynucleotide sequence shown in SEQ ID NOS:1-744, or portion thereof, labeled according to any appropriate method (such as the 32 P-multiprime DNA labeling system (Amersham)) is used as probe. After hybridization overnight, the filter is washed and exposed to x-ray film. DNA for use as probe according to the present invention is described in the sections above and will preferably at least 15 nucleotides in length.

S1 mapping can be performed as described in Fujita et al. (1987) Cell 49:357-367. To prepare probe DNA for use in S1 mapping, the sense strand of an above-described *T. pallidum* DNA sequence of the present invention is used as a template to synthesize labeled antisense DNA. The antisense DNA can then be digested using an appropriate restriction endonuclease to generate further DNA probes of a desired length. Such antisense probes are useful for visualizing protected bands corresponding to the target mRNA (*i.e.*, mRNA encoding *Borrelia* polypeptides).

Levels of mRNA encoding *Borrelia* polypeptides are assayed, for *e.g.*, using the RT-PCR method described in Makino et al. (1990) Technique 2:295-301. By this method, the radioactivities of the "amplicons" in the polyacrylamide gel bands are linearly related to the initial concentration of the target mRNA. Briefly, this method involves adding total RNA isolated from a biological sample in a reaction mixture containing a RT primer and appropriate buffer. After incubating for primer annealing, the mixture can be supplemented with a RT buffer, dNTPs, DTT, RNase inhibitor and reverse transcriptase. After incubation to achieve reverse transcription of the RNA, the RT products are then subject to PCR using labeled primers. Alternatively, rather than labeling the primers, a labeled dNTP can be included in the PCR reaction mixture. PCR amplification can be performed in a DNA thermal cycler according to conventional techniques. After a suitable number of rounds to achieve amplification, the PCR reaction mixture is electrophoresed on a polyacrylamide gel. After drying the gel, the radioactivity of the appropriate bands (corresponding to the mRNA encoding the *Borrelia* polypeptides of the present invention) are quantified using an imaging analyzer. RT and PCR reaction ingredients and conditions, reagent and gel concentrations, and labeling methods are well known in the art. Variations on the RT-PCR method will be apparent to the skilled artisan. Other PCR methods that can detect the nucleic acid of the present invention can be found in PCR PRIMER: A LABORATORY MANUAL (C.W. Dieffenbach et al. eds., Cold Spring Harbor Lab Press, 1995).

The polynucleotides of the present invention, including both DNA and RNA, may be used to detect polynucleotides of the present invention or *Borrelia* species including *T. pallidum* using bio chip technology. The present invention includes both high density chip arrays (>1000 oligonucleotides per cm^2) and low density chip arrays (<1000 oligonucleotides per cm^2). Bio chips comprising arrays of polynucleotides of the present invention may be used to detect *Borrelia* species, including *T. pallidum*, in biological and environmental samples and to diagnose an animal, including humans, with an *T. pallidum* or other *Borrelia* infection. The bio chips of the present invention may comprise polynucleotide sequences of other pathogens including

bacteria, viral, parasitic, and fungal polynucleotide sequences, in addition to the polynucleotide sequences of the present invention, for use in rapid differential pathogenic detection and diagnosis. The bio chips can also be used to monitor an *T. pallidum* or other *Borrelia* infections and to monitor the genetic changes (deletions, insertions, mismatches, etc.) in response to drug therapy in the clinic and drug development in the laboratory. The bio chip technology comprising arrays of polynucleotides of the present invention may also be used to simultaneously monitor the expression of a multiplicity of genes, including those of the present invention. The polynucleotides used to comprise a selected array may be specified in the same manner as for the fragments, i.e., by their 5' and 3' positions or length in contiguous base pairs and include from:

5 Methods and particular uses of the polynucleotides of the present invention to detect *Borrelia* species, including *T. pallidum*, using bio chip technology include those known in the art and those of: U.S. Patent Nos. 5510270, 5545531, 5445934, 5677195, 5532128, 5556752, 10 5527681, 5451683, 5424186, 5607646, 5658732 and World Patent Nos. WO/9710365, WO/9511995, WO/9743447, WO/9535505, each incorporated herein in their entireties.

15 Biosensors using the polynucleotides of the present invention may also be used to detect, diagnose, and monitor *T. pallidum* or other *Borrelia* species and infections thereof. Biosensors using the polynucleotides of the present invention may also be used to detect particular polynucleotides of the present invention. Biosensors using the polynucleotides of the present invention may also be used to monitor the genetic changes (deletions, insertions, mismatches, etc.) in response to drug therapy in the clinic and drug development in the laboratory. Methods and particular uses of the polynucleotides of the present invention to detect *Borrelia* species, including *T. pallidum*, using biosenors include those known in the art and those of: U.S. Patent Nos 5721102, 5658732, 5631170, and World Patent Nos. WO97/35011, WO/9720203, each incorporated herein in their entireties.

20

25 Thus, the present invention includes both bio chips and biosensors comprising polynucleotides of the present invention and methods of their use.

Assaying *Borrelia* polypeptide levels in a biological sample can occur using any art-known method, such as antibody-based techniques. For example, *Borrelia* polypeptide expression in tissues can be studied with classical immunohistological methods. In these, the 30 specific recognition is provided by the primary antibody (polyclonal or monoclonal) but the secondary detection system can utilize fluorescent, enzyme, or other conjugated secondary antibodies. As a result, an immunohistological staining of tissue section for pathological examination is obtained. Tissues can also be extracted, e.g., with urea and neutral detergent, for the liberation of *Borrelia* polypeptides for Western-blot or dot/slot assay. See, e.g., Jalkanen, 35 M. et al. (1985) J. Cell. Biol. 101:976-985; Jalkanen, M. et al. (1987) J. Cell. Biol. 105:3087-3096. In this technique, which is based on the use of cationic solid phases, quantitation of a *Borrelia* polypeptide can be accomplished using an isolated *Borrelia* polypeptide as a standard. This technique can also be applied to body fluids.

Other antibody-based methods useful for detecting *Borrelia* polypeptide gene expression include immunoassays, such as the ELISA and the radioimmunoassay (RIA). For example, a *Borrelia* polypeptide-specific monoclonal antibodies can be used both as an immunoabsorbent and as an enzyme-labeled probe to detect and quantify a *Borrelia* polypeptide. The amount of a 5 *Borrelia* polypeptide present in the sample can be calculated by reference to the amount present in a standard preparation using a linear regression computer algorithm. Such an ELISA is described in Iacobelli et al. (1988) Breast Cancer Research and Treatment 11:19-30. In another ELISA assay, two distinct specific monoclonal antibodies can be used to detect *Borrelia* polypeptides in a body fluid. In this assay, one of the antibodies is used as the immunoabsorbent and the other as 10 the enzyme-labeled probe.

The above techniques may be conducted essentially as a "one-step" or "two-step" assay. The "one-step" assay involves contacting the *Borrelia* polypeptide with immobilized antibody and, without washing, contacting the mixture with the labeled antibody. The "two-step" assay involves washing before contacting the mixture with the labeled antibody. Other conventional 15 methods may also be employed as suitable. It is usually desirable to immobilize one component of the assay system on a support, thereby allowing other components of the system to be brought into contact with the component and readily removed from the sample. Variations of the above and other immunological methods included in the present invention can also be found in Harlow et al., ANTIBODIES: A LABORATORY MANUAL, (Cold Spring Harbor Laboratory Press, 20 2nd ed. 1988).

Suitable enzyme labels include, for example, those from the oxidase group, which catalyze the production of hydrogen peroxide by reacting with substrate. Glucose oxidase is particularly preferred as it has good stability and its substrate (glucose) is readily available. Activity of an oxidase label may be assayed by measuring the concentration of hydrogen peroxide 25 formed by the enzyme-labeled antibody/substrate reaction. Besides enzymes, other suitable labels include radioisotopes, such as iodine (^{125}I , ^{131}I), carbon (^{14}C), sulphur (^{35}S), tritium (^3H), indium (^{113}In), and technetium (^{99m}Tc), and fluorescent labels, such as fluorescein and rhodamine, and biotin.

Further suitable labels for the *Borrelia* polypeptide-specific antibodies of the present 30 invention are provided below. Examples of suitable enzyme labels include malate dehydrogenase, *Borrelia* nuclease, delta-5-steroid isomerase, yeast-alcohol dehydrogenase, alpha-glycerol phosphate dehydrogenase, triose phosphate isomerase, peroxidase, alkaline phosphatase, asparaginase, glucose oxidase, beta-galactosidase, ribonuclease, urease, catalase, glucose-6-phosphate dehydrogenase, glucoamylase, and acetylcholine esterase.

Examples of suitable radioisotopic labels include ^3H , ^{111}In , ^{125}I , ^{131}I , ^{32}P , ^{35}S , ^{14}C , ^{51}Cr , 35 ^{57}Co , ^{58}Co , ^{59}Fe , ^{75}Se , ^{152}Eu , ^{90}Y , ^{67}Cu , ^{217}Bi , ^{211}At , ^{212}Pb , ^{47}Sc , ^{109}Pd , etc. ^{111}In is a preferred isotope where *in vivo* imaging is used since its avoids the problem of dehalogenation of the ^{125}I or ^{131}I -labeled monoclonal antibody by the liver. In addition, this radionucleotide has a more favorable gamma emission energy for imaging. See, e.g., Perkins et al. (1985) Eur. J. Nucl.

Med. 10:296-301; Carasquillo et al. (1987) J. Nucl. Med. 28:281-287. For example, ¹¹¹In coupled to monoclonal antibodies with 1-(P-isothiocyanatobenzyl)-DPTA has shown little uptake in non-tumors tissues, particularly the liver, and therefore enhances specificity of tumor localization. See, Esteban et al. (1987) J. Nucl. Med. 28:861-870.

5 Examples of suitable non-radioactive isotopic labels include ¹⁵⁷Gd, ⁵⁵Mn, ¹⁶²Dy, ⁵²Tr, and ⁵⁶Fe.

Examples of suitable fluorescent labels include an ¹⁵²Eu label, a fluorescein label, an isothiocyanate label, a rhodamine label, a phycoerythrin label, a phycocyanin label, an allophycocyanin label, an o-phthaldehyde label, and a fluorescamine label.

10 Examples of suitable toxin labels include, *Pseudomonas* toxin, diphtheria toxin, ricin, and cholera toxin.

Examples of chemiluminescent labels include a luminal label, an isoluminal label, an aromatic acridinium ester label, an imidazole label, an acridinium salt label, an oxalate ester label, a luciferin label, a luciferase label, and an aequorin label.

15 Examples of nuclear magnetic resonance contrasting agents include heavy metal nuclei such as Gd, Mn, and iron.

Typical techniques for binding the above-described labels to antibodies are provided by Kennedy et al. (1976) Clin. Chim. Acta 70:1-31, and Schurs et al. (1977) Clin. Chim. Acta 81:1-40. Coupling techniques mentioned in the latter are the glutaraldehyde method, the 20 periodate method, the dimaleimide method, the m-maleimidobenzyl-N-hydroxy-succinimide ester method, all of which methods are incorporated by reference herein.

25 In a related aspect, the invention includes a diagnostic kit for use in screening serum containing antibodies specific against *T. pallidum* infection. Such a kit may include an isolated *T. pallidum* antigen comprising an epitope which is specifically immunoreactive with at least one anti-*T. pallidum* antibody. Such a kit also includes means for detecting the binding of said antibody to the antigen. In specific embodiments, the kit may include a recombinantly produced or chemically synthesized peptide or polypeptide antigen. The peptide or polypeptide antigen may be attached to a solid support.

30 In a more specific embodiment, the detecting means of the above-described kit includes a solid support to which said peptide or polypeptide antigen is attached. Such a kit may also include a non-attached reporter-labeled anti-human antibody. In this embodiment, binding of the antibody to the *T. pallidum* antigen can be detected by binding of the reporter labeled antibody to the anti-*T. pallidum* polypeptide antibody.

35 Specifically, the invention provides a compartmentalized kit to receive, in close confinement, one or more containers which comprises: (a) a first container comprising one of the DFs or antibodies of the present invention; and (b) one or more other containers comprising one or more of the following: wash reagents, reagents capable of detecting presence of a bound DF or antibody.

In detail, a compartmentalized kit includes any kit in which reagents are contained in separate containers. Such containers include small glass containers, plastic containers or strips of plastic or paper. Such containers allows one to efficiently transfer reagents from one compartment to another compartment such that the samples and reagents are not cross-

5 contaminated, and the agents or solutions of each container can be added in a quantitative fashion from one compartment to another. Such containers will include a container which will accept the test sample, a container which contains the antibodies used in the assay, containers which contain wash reagents (such as phosphate buffered saline, Tris-buffers, etc.), and containers which contain the reagents used to detect the bound antibody or DF.

10

In a related aspect, the invention includes a method of detecting *T. pallidum* infection in a subject. This detection method includes reacting a body fluid, preferably serum, from the subject with an isolated *T. pallidum* antigen, and examining the antigen for the presence of bound antibody. In a specific embodiment, the method includes a polypeptide antigen attached to a solid support, and serum is reacted with the support. Subsequently, the support is reacted with a reporter-labeled anti-human antibody. The support is then examined for the presence of reporter-labeled antibody.

15 The solid surface reagent employed in the above assays and kits is prepared by known techniques for attaching protein material to solid support material, such as polymeric beads, dip sticks, 96-well plates or filter material. These attachment methods generally include non-specific adsorption of the protein to the support or covalent attachment of the protein, typically through a free amine group, to a chemically reactive group on the solid support, such as an activated carboxyl, hydroxyl, or aldehyde group. Alternatively, streptavidin coated plates can be used in conjunction with biotinylated antigen(s).

20 The polypeptides and antibodies of the present invention, including fragments thereof, may be used to detect *Borrelia* species including *T. pallidum* using bio chip and biosensor technology. Bio chip and biosensors of the present invention may comprise the polypeptides of the present invention to detect antibodies, which specifically recognize *Borrelia* species, including *T. pallidum*. Bio chip and biosensors of the present invention may also comprise antibodies 25 which specifically recognize the polypeptides of the present invention to detect *Borrelia* species, including *T. pallidum* or specific polypeptides of the present invention. Bio chips or biosensors comprising polypeptides or antibodies of the present invention may be used to detect *Borrelia* species, including *T. pallidum*, in biological and environmental samples and to diagnose an animal, including humans, with an *T. pallidum* or other *Borrelia* infection. Thus, the present 30 invention includes both bio chips and biosensors comprising polypeptides or antibodies of the present invention and methods of their use.

35 The bio chips of the present invention may further comprise polypeptide sequences of other pathogens including bacteria, viral, parasitic, and fungal polypeptide sequences, in addition to the polypeptide sequences of the present invention, for use in rapid differential pathogenic

detection and diagnosis. The bio chips of the present invention may further comprise antibodies or fragments thereof specific for other pathogens including bacteria, viral, parasitic, and fungal polypeptide sequences, in addition to the antibodies or fragments thereof of the present invention, for use in rapid differential pathogenic detection and diagnosis. The bio chips and biosensors of the present invention may also be used to monitor an *T. pallidum* or other *Borrelia* infection and to monitor the genetic changes (amino acid deletions, insertions, substitutions, etc.) in response to drug therapy in the clinic and drug development in the laboratory. The bio chip and biosensors comprising polypeptides or antibodies of the present invention may also be used to simultaneously monitor the expression of a multiplicity of polypeptides, including those of the present invention. The polypeptides used to comprise a bio chip or biosensor of the present invention may be specified in the same manner as for the fragments, i.e., by their N-terminal and C-terminal positions or length in contiguous amino acid residue. Methods and particular uses of the polypeptides and antibodies of the present invention to detect *Borrelia* species, including *T. pallidum*, or specific polypeptides using bio chip and biosensor technology include those known in the art, those of the U.S. Patent Nos. and World Patent Nos. listed above for bio chips and biosensors using polynucleotides of the present invention, and those of: U.S. Patent Nos. 5658732, 5135852, 5567301, 5677196, 5690894 and World Patent Nos. WO9729366, WO9612957, each incorporated herein in their entireties.

20 **4. Screening Assay for Binding Agents**

Using the isolated proteins of the present invention, the present invention further provides methods of obtaining and identifying agents which bind to a protein encoded by one of the ORFs of the present invention or to one of the fragments and the *T. pallidum* fragment and contigs herein described.

25 In general, such methods comprise steps of:

- (a) contacting an agent with an isolated protein encoded by one of the ORFs of the present invention, or an isolated fragment of the *T. pallidum* genome; and
- (b) determining whether the agent binds to said protein or said fragment.

The agents screened in the above assay can be, but are not limited to, peptides, carbohydrates, vitamin derivatives, or other pharmaceutical agents. The agents can be selected and screened at random or rationally selected or designed using protein modeling techniques.

For random screening, agents such as peptides, carbohydrates, pharmaceutical agents and the like are selected at random and are assayed for their ability to bind to the protein encoded by the ORF of the present invention.

30 Alternatively, agents may be rationally selected or designed. As used herein, an agent is said to be "rationally selected or designed" when the agent is chosen based on the configuration of the particular protein. For example, one skilled in the art can readily adapt currently available procedures to generate peptides, pharmaceutical agents and the like capable of binding to a specific peptide sequence in order to generate rationally designed antipeptide peptides, for

example see Hurby *et al.*, "Application of Synthetic Peptides: Antisense Peptides," in *Synthetic Peptides, A User's Guide*, W. H. Freeman, NY (1992), pp. 289-307, and Kaspaczak *et al.*, *Biochemistry* 28:9230-8 (1989), or pharmaceutical agents, or the like.

In addition to the foregoing, one class of agents of the present invention, as broadly described, can be used to control gene expression through binding to one of the ORFs or EMFs of the present invention. As described above, such agents can be randomly screened or rationally designed/selected. Targeting the ORF or EMF allows a skilled artisan to design sequence specific or element specific agents, modulating the expression of either a single ORF or multiple ORFs which rely on the same EMF for expression control.

One class of DNA binding agents are agents which contain base residues which hybridize or form a triple helix by binding to DNA or RNA. Such agents can be based on the classic phosphodiester, ribonucleic acid backbone, or can be a variety of sulphydryl or polymeric derivatives which have base attachment capacity.

Agents suitable for use in these methods usually contain 20 to 40 bases and are designed to be complementary to a region of the gene involved in transcription (triple helix - see Lee *et al.*, *Nucl. Acids Res.* 6:3073 (1979); Cooney *et al.*, *Science* 241:456 (1988); and Dervan *et al.*, *Science* 251:1360 (1991)) or to the mRNA itself (antisense - Okano, *J. Neurochem.* 56:560 (1991); *Oligodeoxynucleotides as Antisense Inhibitors of Gene Expression*, CRC Press, Boca Raton, FL (1988)). Triple helix- formation optimally results in a shut-off of RNA transcription from DNA, while antisense RNA hybridization blocks translation of an mRNA molecule into polypeptide. Both techniques have been demonstrated to be effective in model systems. Information contained in the sequences of the present invention can be used to design antisense and triple helix-forming oligonucleotides, and other DNA binding agents.

5. Pharmaceutical Compositions and Vaccines

The present invention further provides pharmaceutical agents which can be used to modulate the growth or pathogenicity of *T. pallidum*, or another related organism, *in vivo* or *in vitro*. As used herein, a "pharmaceutical agent" is defined as a composition of matter which can be formulated using known techniques to provide a pharmaceutical compositions. As used herein, the "pharmaceutical agents of the present invention" refers the pharmaceutical agents which are derived from the proteins encoded by the ORFs of the present invention or are agents which are identified using the herein described assays.

As used herein, a pharmaceutical agent is said to "modulate the growth pathogenicity of *T. pallidum* or a related organism, *in vivo* or *in vitro*," when the agent reduces the rate of growth, rate of division, or viability of the organism in question. The pharmaceutical agents of the present invention can modulate the growth or pathogenicity of an organism in many fashions, although an understanding of the underlying mechanism of action is not needed to practice the use of the pharmaceutical agents of the present invention. Some agents will modulate the growth by binding to an important protein thus blocking the biological activity of the protein, while other

agents may bind to a component of the outer surface of the organism blocking attachment or rendering the organism more prone to act the bodies nature immune system. Alternatively, the agent may comprise a protein encoded by one of the ORFs of the present invention and serve as a vaccine. The development and use of a vaccine based on outer membrane components are well known in the art.

As used herein, a "related organism" is a broad term which refers to any organism whose growth can be modulated by one of the pharmaceutical agents of the present invention. In general, such an organism will contain a homolog of the protein which is the target of the pharmaceutical agent or the protein used as a vaccine. As such, related organisms do not need to be bacterial but may be fungal or viral pathogens.

The pharmaceutical agents and compositions of the present invention may be administered in a convenient manner, such as by the oral, topical, intravenous, intraperitoneal, intramuscular, subcutaneous, intranasal or intradermal routes. The pharmaceutical compositions are administered in an amount which is effective for treating and/or prophylaxis of the specific indication. In general, they are administered in an amount of at least about 1 mg/kg body weight and in most cases they will be administered in an amount not in excess of about 1 g/kg body weight per day. In most cases, the dosage is from about 0.1 mg/kg to about 10 g/kg body weight daily, taking into account the routes of administration, symptoms, etc.

The agents of the present invention can be used in native form or can be modified to form a chemical derivative. As used herein, a molecule is said to be a "chemical derivative" of another molecule when it contains additional chemical moieties not normally a part of the molecule. Such moieties may improve the molecule's solubility, absorption, biological half life, etc. The moieties may alternatively decrease the toxicity of the molecule, eliminate or attenuate any undesirable side effect of the molecule, etc. Moieties capable of mediating such effects are disclosed in, among other sources, REMINGTON'S PHARMACEUTICAL SCIENCES (1980) cited elsewhere herein.

For example, such moieties may change an immunological character of the functional derivative, such as affinity for a given antibody. Such changes in immunomodulation activity are measured by the appropriate assay, such as a competitive type immunoassay. Modifications of such protein properties as redox or thermal stability, biological half-life, hydrophobicity, susceptibility to proteolytic degradation or the tendency to aggregate with carriers or into multimers also may be effected in this way and can be assayed by methods well known to the skilled artisan.

The therapeutic effects of the agents of the present invention may be obtained by providing the agent to a patient by any suitable means (e.g., inhalation, intravenously, intramuscularly, subcutaneously, enterally, or parenterally). It is preferred to administer the agent of the present invention so as to achieve an effective concentration within the blood or tissue in which the growth of the organism is to be controlled. To achieve an effective blood

concentration, the preferred method is to administer the agent by injection. The administration may be by continuous infusion, or by single or multiple injections.

In providing a patient with one of the agents of the present invention, the dosage of the administered agent will vary depending upon such factors as the patient's age, weight, height, sex, general medical condition, previous medical history, etc. In general, it is desirable to provide the recipient with a dosage of agent which is in the range of from about 1 pg/kg to 10 mg/kg (body weight of patient), although a lower or higher dosage may be administered. The therapeutically effective dose can be lowered by using combinations of the agents of the present invention or another agent.

As used herein, two or more compounds or agents are said to be administered "in combination" with each other when either (1) the physiological effects of each compound, or (2) the serum concentrations of each compound can be measured at the same time. The composition of the present invention can be administered concurrently with, prior to, or following the administration of the other agent.

The agents of the present invention are intended to be provided to recipient subjects in an amount sufficient to decrease the rate of growth (as defined above) of the target organism.

The administration of the agent(s) of the invention may be for either a "prophylactic" or "therapeutic" purpose. When provided prophylactically, the agent(s) are provided in advance of any symptoms indicative of the organisms growth. The prophylactic administration of the

agent(s) serves to prevent, attenuate, or decrease the rate of onset of any subsequent infection. When provided therapeutically, the agent(s) are provided at (or shortly after) the onset of an indication of infection. The therapeutic administration of the compound(s) serves to attenuate the pathological symptoms of the infection and to increase the rate of recovery.

The agents of the present invention are administered to a subject, such as a mammal, or a patient, in a pharmaceutically acceptable form and in a therapeutically effective concentration. A composition is said to be "pharmacologically acceptable" if its administration can be tolerated by a recipient patient. Such an agent is said to be administered in a "therapeutically effective amount" if the amount administered is physiologically significant. An agent is physiologically significant if its presence results in a detectable change in the physiology of a recipient patient.

The agents of the present invention can be formulated according to known methods to prepare pharmaceutically useful compositions, whereby these materials, or their functional derivatives, are combined in a mixture with a pharmaceutically acceptable carrier vehicle. Suitable vehicles and their formulation, inclusive of other human proteins, e.g., human serum albumin, are described, for example, in REMINGTON'S PHARMACEUTICAL SCIENCES, 16th Ed., Osol, A., Ed., Mack Publishing, Easton PA (1980). In order to form a pharmaceutically acceptable composition suitable for effective administration, such compositions will contain an effective amount of one or more of the agents of the present invention, together with a suitable amount of carrier vehicle.

Additional pharmaceutical methods may be employed to control the duration of action. Control release preparations may be achieved through the use of polymers to complex or absorb one or more of the agents of the present invention. The controlled delivery may be effectuated by a variety of well known techniques, including formulation with macromolecules such as, for example, polyesters, polyamino acids, polyvinyl, pyrrolidone, ethylenevinylacetate, methylcellulose, carboxymethylcellulose, or protamine, sulfate, adjusting the concentration of the macromolecules and the agent in the formulation, and by appropriate use of methods of incorporation, which can be manipulated to effectuate a desired time course of release. Another possible method to control the duration of action by controlled release preparations is to incorporate agents of the present invention into particles of a polymeric material such as polyesters, polyamino acids, hydrogels, poly(lactic acid) or ethylene vinylacetate copolymers. Alternatively, instead of incorporating these agents into polymeric particles, it is possible to entrap these materials in microcapsules prepared, for example, by coacervation techniques or by interfacial polymerization with, for example, hydroxymethylcellulose or gelatine-microcapsules and poly(methylmethacrylate) microcapsules, respectively, or in colloidal drug delivery systems, for example, liposomes, albumin microspheres, microemulsions, nanoparticles, and nanocapsules or in macroemulsions. Such techniques are disclosed in REMINGTON'S PHARMACEUTICAL SCIENCES (1980).

The invention further provides a pharmaceutical pack or kit comprising one or more containers filled with one or more of the ingredients of the pharmaceutical compositions of the invention. Associated with such container(s) can be a notice in the form prescribed by a governmental agency regulating the manufacture, use or sale of pharmaceuticals or biological products, which notice reflects approval by the agency of manufacture, use or sale for human administration.

In addition, the agents of the present invention may be employed in conjunction with other therapeutic compounds.

6. Shot-Gun Approach to Megabase DNA Sequencing

The present invention further demonstrates that a large sequence can be sequenced using a random shotgun approach. This procedure, described in detail in the examples that follow, has eliminated the up front cost of isolating and ordering overlapping or contiguous subclones prior to the start of the sequencing protocols.

Certain aspects of the present invention are described in greater detail in the examples that follow. The examples are provided by way of illustration. Other aspects and embodiments of the present invention are contemplated by the inventors, as will be clear to those of skill in the art from reading the present disclosure.

ILLUSTRATIVE EXAMPLES

LIBRARIES AND SEQUENCING

1. Shotgun Sequencing Probability Analysis

The overall strategy for a shotgun approach to whole genome sequencing follows from

5 the Lander and Waterman (Landerman and Waterman, *Genomics* 2:231 (1988)) application of the equation for the Poisson distribution. According to this treatment, the probability, P0, that any given base in a sequence of size L, in nucleotides, is not sequenced after a certain amount, n, in nucleotides, of random sequence has been determined can be calculated by the equation $P0 = e^{-m}$, where m is L/n , the fold coverage. For instance, for a genome of 2.8 Mb, m=1 when 2.8

10 Mb of sequence has been randomly generated (1X coverage). At that point, $P0 = e^{-1} = 0.37$.

The probability that any given base has not been sequenced is the same as the probability that any region of the whole sequence L has not been determined and, therefore, is equivalent to the fraction of the whole sequence that has yet to be determined. Thus, at one-fold coverage, approximately 37% of a polynucleotide of size L, in nucleotides has not been sequenced. When

15 14 Mb of sequence has been generated, coverage is 5X for a 2.8 Mb and the unsequenced fraction drops to .0067 or 0.67%. 5X coverage of a 2.8 Mb sequence can be attained by sequencing approximately 17,000 random clones from both insert ends with an average sequence read length of 410 bp.

Similarly, the total gap length, G, is determined by the equation $G = Le-m$, and the

20 average gap size, g, follows the equation, $g = L/n$. Thus, 5X coverage leaves about 240 gaps averaging about 82 bp in size in a sequence of a polynucleotide 2.8 Mb long.

The treatment above is essentially that of Lander and Waterman, *Genomics* 2: 231 (1988).

25 2. Random Library Construction

In order to approximate the random model described above during actual sequencing, a nearly ideal library of cloned genomic fragments is required. The following library construction procedure was developed to achieve this end.

T. pallidum DNA is prepared by phenol extraction. A mixture containing 200 µg DNA in

30 1.0 ml of 300 mM sodium acetate, 10 mM Tris-HCl, 1 mM Na-EDTA, 50% glycerol is processed through a nebulizer (IPI Medical Products) with a stream of nitrogen adjusted to 35 Kpa for 2 minutes. The sonicated DNA is ethanol precipitated and redissolved in 500 µl TE buffer.

To create blunt-ends, a 100 µl aliquot of the resuspended DNA is digested with 5 units of

35 BAL31 nuclease (New England BioLabs) for 10 min at 30°C in 200 µl BAL31 buffer. The digested DNA is phenol-extracted, ethanol-precipitated, redissolved in 100 µl TE buffer, and then size-fractionated by electrophoresis through a 1.0% low melting temperature agarose gel. The section containing DNA fragments 1.6-2.0 kb in size is excised from the gel, and the LGT agarose is melted and the resulting solution is extracted with phenol to separate the agarose from

the DNA. DNA is ethanol precipitated and redissolved in 20 μ l of TE buffer for ligation to vector.

A two-step ligation procedure is used to produce a plasmid library with 97% inserts, of which >99% were single inserts. The first ligation mixture (50 μ l) contains 2 μ g of DNA

- 5 fragments, 2 μ g pUC18 DNA (Pharmacia) cut with SmaI and dephosphorylated with bacterial alkaline phosphatase, and 10 units of T4 ligase (GIBCO/BRL) and is incubated at 14°C for 4 hr. The ligation mixture then is phenol extracted and ethanol precipitated, and the precipitated DNA is dissolved in 20 μ l TE buffer and electrophoresed on a 1.0% low melting agarose gel. Discrete bands in a ladder are visualized by ethidium bromide-staining and UV illumination and identified
10 by size as insert (I), vector (v), v+I, v+2i, v+3i, etc. The portion of the gel containing v+I DNA is excised and the v+I DNA is recovered and resuspended into 20 μ l TE. The v+I DNA then is blunt-ended by T4 polymerase treatment for 5 min. at 37°C in a reaction mixture (50 μ l) containing the v+I linear, 500 μ M each of the 4 dNTPs, and 9 units of T4 polymerase (New England BioLabs), under recommended buffer conditions. After phenol extraction and ethanol
15 precipitation the repaired v+I linear are dissolved in 20 μ l TE. The final ligation to produce circles is carried out in a 50 μ l reaction containing 5 μ l of v+I linear and 5 units of T4 ligase at 14°C overnight. After 10 min. at 70°C the following day, the reaction mixture is stored at -20°C.

This two-stage procedure results in a molecularly random collection of single-insert plasmid recombinants with minimal contamination from double-insert chimeras (<1%) or free
20 vector (<3%).

Since deviation from randomness can arise from propagation the DNA in the host, *E. coli* host cells deficient in all recombination and restriction functions (A. Greener, *Strategies* 3 (1):5 (1990)) are used to prevent rearrangements, deletions, and loss of clones by restriction. Furthermore, transformed cells are plated directly on antibiotic diffusion plates to avoid the usual
25 broth recovery phase which allows multiplication and selection of the most rapidly growing cells.

Plating is carried out as follows. A 100 μ l aliquot of Epicurian Coli SURE II Supercompetent Cells (Stratagene 200152) is thawed on ice and transferred to a chilled Falcon 2059 tube on ice. A 1.7 μ l aliquot of 1.42 M beta-mercaptoethanol is added to the aliquot of cells to a final concentration of 25 mM. Cells are incubated on ice for 10 min. A 1 μ l aliquot of the final ligation is added to the cells and incubated on ice for 30 min. The cells are heat pulsed for 30 sec. at 42°C and placed back on ice for 2 min. The outgrowth period in liquid culture is eliminated from this protocol in order to minimize the preferential growth of any given transformed cell. Instead the transformation mixture is plated directly on a nutrient rich SOB plate containing a 5 ml bottom layer of SOB agar (5% SOB agar: 20 g tryptone, 5 g yeast extract, 30 0.5 g NaCl, 1.5% Difco Agar per liter of media). The 5 ml bottom layer is supplemented with 0.4 ml of 50 mg/ml ampicillin per 100 ml SOB agar. The 15 ml top layer of SOB agar is supplemented with 1 ml X-Gal (2%), 1 ml MgCl₂ (1 M), and 1 ml MgSO₄/100 ml SOB agar. The 15 ml top layer is poured just prior to plating. Our titer is approximately 100 colonies/10 μ l aliquot of transformation.

All colonies are picked for template preparation regardless of size. Thus, only clones lost due to "poison" DNA or deleterious gene products are deleted from the library, resulting in a slight increase in gap number over that expected.

5 3. Random DNA Sequencing

High quality double stranded DNA plasmid templates are prepared using a "boiling bead" method developed in collaboration with Advanced Genetic Technology Corp. (Gaithersburg, MD) (Adams *et al.*, *Science* 252:1651 (1991); Adams *et al.*, *Nature* 355:632 (1992)). Plasmid preparation is performed in a 96-well format for all stages of DNA preparation from bacterial growth through final DNA purification. Template concentration is determined using Hoechst Dye and a Millipore Cytofluor. DNA concentrations are not adjusted, but low-yielding templates are identified where possible and not sequenced.

Templates are also prepared from two *T. pallidum* lambda genomic libraries. An amplified library is constructed in the vector Lambda GEM-12 (Promega) and an unamplified library is constructed in Lambda DASH II (Stratagene). In particular, for the unamplified lambda library, *T. pallidum* DNA (> 100 kb) is partially digested in a reaction mixture (200 μ l) containing 50 μ g DNA, 1X Sau3AI buffer, 20 units Sau3AI for 6 min. at 23°C. The digested DNA was phenol-extracted and electrophoresed on a 0.5% low melting agarose gel at 2V/cm for 7 hours. Fragments from 15 to 25 kb are excised and recovered in a final volume of 6 μ l. One μ l of fragments is used with 1 μ l of DASHII vector (Stratagene) in the recommended ligation reaction. One μ l of the ligation mixture is used per packaging reaction following the recommended protocol with the Gigapack II XL Packaging Extract (Stratagene, #227711). Phage are plated directly without amplification from the packaging mixture (after dilution with 500 μ l of recommended SM buffer and chloroform treatment). Yield is about 2.5x10³ pfu/ μ l. The amplified library is prepared essentially as above except the lambda GEM-12 vector is used. After packaging, about 3.5x10⁴ pfu are plated on the restrictive NM539 host. The lysate is harvested in 2 ml of SM buffer and stored frozen in 7% dimethylsulfoxide. The phage titer is approximately 1x10⁹ pfu/ml.

Liquid lysates (100 μ l) are prepared from randomly selected plaques (from the unamplified library) and template is prepared by long-range PCR using T7 and T3 vector-specific primers.

Sequencing reactions are carried out on plasmid and/or PCR templates using the AB Catalyst LabStation with Applied Biosystems PRISM Ready Reaction Dye Primer Cycle Sequencing Kits for the M13 forward (M13-21) and the M13 reverse (M13RP1) primers (Adams *et al.*, *Nature* 368:474 (1994)). Dye terminator sequencing reactions are carried out on the lambda templates on a Perkin-Elmer 9600 Thermocycler using the Applied Biosystems Ready Reaction Dye Terminator Cycle Sequencing kits. T7 and SP6 primers are used to sequence the ends of the inserts from the Lambda GEM-12 library and T7 and T3 primers are used to sequence the ends of the inserts from the Lambda DASH II library. Sequencing reactions are performed

by eight individuals using an average of fourteen AB 373 DNA Sequencers per day. All sequencing reactions are analyzed using the Stretch modification of the AB 373, primarily using a 34 cm well-to-read distance. The overall sequencing success rate very approximately is about 85% for M13-21 and M13RP1 sequences and 65% for dye-terminator reactions. The average 5 usable read length is 485 bp for M13-21 sequences, 445bp for M13RP1 sequences, and 375 bp for dye-terminator reactions.

Richards *et al.*, Chapter 28 in AUTOMATED DNA SEQUENCING AND ANALYSIS, M. D. Adams, C. Fields, J. C. Venter, Eds., Academic Press, London, (1994) described the value of using sequence from both ends of sequencing templates to facilitate ordering of contigs 10 in shotgun assembly projects of lambda and cosmid clones. We balance the desirability of both-end sequencing (including the reduced cost of lower total number of templates) against shorter read-lengths for sequencing reactions performed with the M13RP1 (reverse) primer compared to the M13-21 (forward) primer. Approximately one-half of the templates are sequenced from both 15 ends. Random reverse sequencing reactions are done based on successful forward sequencing reactions. Some M13RP1 sequences are obtained in a semi-directed fashion: M13-21: sequences pointing outward at the ends of contigs are chosen for M13RP1 sequencing in an effort to specifically order contigs.

4. Protocol for Automated Cycle Sequencing

20 The sequencing is carried out using ABI Catalyst robots and AB 373 Automated DNA Sequencers. The Catalyst robot is a publicly available sophisticated pipetting and temperature control robot which has been developed specifically for DNA sequencing reactions. The Catalyst combines pre-aliquoted templates and reaction mixes consisting of deoxy- and dideoxynucleotides, the thermostable Taq DNA polymerase, fluorescently-labelled sequencing 25 primers, and reaction buffer. Reaction mixes and templates are combined in the wells of an aluminum 96-well thermocycling plate. Thirty consecutive cycles of linear amplification (*i.e.*., one primer synthesis) steps are performed including denaturation, annealing of primer and template, and extension; *i.e.*, DNA synthesis. A heated lid with rubber gaskets on the thermocycling plate prevents evaporation without the need for an oil overlay.

30 Two sequencing protocols are used: one for dye-labelled primers and a second for dye-labelled dideoxy chain terminators. The shotgun sequencing involves use of four dye-labelled sequencing primers, one for each of the four terminator nucleotide. Each dye-primer is labelled with a different fluorescent dye, permitting the four individual reactions to be combined into one lane of the 373 DNA Sequencer for electrophoresis, detection, and base-calling. ABI currently 35 supplies pre-mixed reaction mixes in bulk packages containing all the necessary non-template reagents for sequencing. Sequencing can be done with both plasmid and PCR- generated templates with both dye-primers and dye- terminators with approximately equal fidelity, although plasmid templates generally give longer usable sequences.

Thirty-two reactions are loaded per AB373 Sequencer each day, for a total of 960 samples. Electrophoresis is run overnight following the manufacturer's protocols, and the data is collected for twelve hours. Following electrophoresis and fluorescence detection, the ABI 373 performs automatic lane tracking and base-calling. The lane-tracking is confirmed visually. Each 5 sequence electropherogram (or fluorescence lane trace) is inspected visually and assessed for quality. Trailing sequences of low quality are removed and the sequence itself is loaded via software to a Sybase database (archived daily to 8mm tape). Leading vector polylinker sequence is removed automatically by a software program. Average edited lengths of sequences from the standard ABI 373 are around 400 bp and depend mostly on the quality of the template used for 10 the sequencing reaction. ABI 373 Sequencers converted to Stretch Liners provide a longer electrophoresis path prior to fluorescence detection and increase the average number of usable bases to 500-600 bp.

INFORMATICS

1. Data Management

A number of information management systems for a large-scale sequencing lab have been developed. (For review see, for instance, Kerlavage *et al.*, *Proceedings of the Twenty-Sixth Annual Hawaii International Conference on System Sciences*, IEEE Computer Society Press, Washington D. C., 585 (1993)) The system used to collect and assemble the sequence data was 20 developed using the Sybase relational database management system and was designed to automate data flow wherever possible and to reduce user error. The database stores and correlates all information collected during the entire operation from template preparation to final analysis of the genome. Because the raw output of the ABI 373 Sequencers was based on a Macintosh platform and the data management system chosen was based on a Unix platform, it 25 was necessary to design and implement a variety of multi-user, client-server applications which allow the raw data as well as analysis results to flow seamlessly into the database with a minimum of user effort.

2. Assembly

An assembly engine (TIGR Assembler) developed for the rapid and accurate assembly of 30 thousands of sequence fragments was employed to generate contigs. The TIGR assembler simultaneously clusters and assembles fragments of the genome. In order to obtain the speed necessary to assemble more than 104 fragments, the algorithm builds a hash table of 12 bp oligonucleotide subsequences to generate a list of potential sequence fragment overlaps. The 35 number of potential overlaps for each fragment determines which fragments are likely to fall into repetitive elements. Beginning with a single seed sequence fragment, TIGR Assembler extends the current contig by attempting to add the best matching fragment based on oligonucleotide content. The contig and candidate fragment are aligned using a modified version of the Smith-Waterman algorithm which provides for optimal gapped alignments (Waterman, M. S., *Methods*

in *Enzymology* 164:765 (1988)). The contig is extended by the fragment only if strict criteria for the quality of the match are met. The match criteria include the minimum length of overlap, the maximum length of an unmatched end, and the minimum percentage match. These criteria are automatically lowered by the algorithm in regions of minimal coverage and raised in regions with 5 a possible repetitive element. The number of potential overlaps for each fragment determines which fragments are likely to fall into repetitive elements. Fragments representing the boundaries of repetitive elements and potentially chimeric fragments are often rejected based on partial mismatches at the ends of alignments and excluded from the current contig. TIGR Assembler is designed to take advantage of clone size information coupled with sequencing from both ends of 10 each template. It enforces the constraint that sequence fragments from two ends of the same template point toward one another in the contig and are located within a certain range of base pairs (definable for each clone based on the known clone size range for a given library).

The process resulted in 744 contigs as represented by SEQ ID NOs:1-744.

15 **3. Identifying Genes**

The predicted coding regions of the *T. pallidum* genome were initially defined with the program GeneMark, which finds ORFs using a probabilistic classification technique. The predicted coding region sequences were used in searches against a database of all nucleotide sequences from GenBank (June, 1997), using the BLASTN search method to identify overlaps 20 of 50 or more nucleotides with at least a 95% identity. Those ORFs with nucleotide sequence matches are shown in Table 1. The ORFs without such matches were translated to protein sequences and compared to a non-redundant database of known proteins generated by combining the Swiss-prot, PIR and GenPept databases. ORFs that matched a database protein with BLASTP probability less than or equal to 0.01 are shown in Table 2. The table also lists 25 assigned functions based on the closest match in the databases. ORFs that did not match protein or nucleotide sequences in the databases at these levels are shown in Table 3.

ILLUSTRATIVE APPLICATIONS

1. Production of an Antibody to a *T. pallidum* Protein

30 Substantially pure protein or polypeptide is isolated from the transfected or transformed cells using any one of the methods known in the art. The protein can also be produced in a recombinant prokaryotic expression system, such as *E. coli*, or can be chemically synthesized. Concentration of protein in the final preparation is adjusted, for example, by concentration on an Amicon filter device, to the level of a few micrograms/ml. Monoclonal or polyclonal antibody to 35 the protein can then be prepared as follows.

2. Monoclonal Antibody Production by Hybridoma Fusion

Monoclonal antibody to epitopes of any of the peptides identified and isolated as described can be prepared from murine hybridomas according to the classical method of Kohler,

G. and Milstein, C., *Nature* 256:495 (1975) or modifications of the methods thereof. Briefly, a mouse is repetitively inoculated with a few micrograms of the selected protein over a period of a few weeks. The mouse is then sacrificed, and the antibody producing cells of the spleen isolated. The spleen cells are fused by means of polyethylene glycol with mouse myeloma cells, 5 and the excess unfused cells destroyed by growth of the system on selective media comprising aminopterin (HAT media). The successfully fused cells are diluted and aliquots of the dilution placed in wells of a microtiter plate where growth of the culture is continued. Antibody-producing clones are identified by detection of antibody in the supernatant fluid of the wells by immunoassay procedures, such as ELISA, as originally described by Engvall, E., *Meth. Enzymol.* 70:419 (1980), and modified methods thereof. Selected positive clones can be 10 expanded and their monoclonal antibody product harvested for use. Detailed procedures for monoclonal antibody production are described in Davis, L. et al., *Basic Methods in Molecular Biology*, Elsevier, New York. Section 21-2 (1989).

15 **3. Polyclonal Antibody Production by Immunization**

Polyclonal antiserum containing antibodies to heterogenous epitopes of a single protein can be prepared by immunizing suitable animals with the expressed protein described above, which can be unmodified or modified to enhance immunogenicity. Effective polyclonal antibody production is affected by many factors related both to the antigen and the host species. For 20 example, small molecules tend to be less immunogenic than others and may require the use of carriers and adjuvant. Also, host animals vary in response to site of inoculations and dose, with both inadequate or excessive doses of antigen resulting in low titer antisera. Small doses (ng level) of antigen administered at multiple intradermal sites appears to be most reliable. An effective immunization protocol for rabbits can be found in Vaitukaitis, J. et al., *J. Clin. Endocrinol. Metab.* 33:988-991 (1971).

25 Booster injections can be given at regular intervals, and antiserum harvested when antibody titer thereof, as determined semi-quantitatively, for example, by double immunodiffusion in agar against known concentrations of the antigen, begins to fall. See, for example, Ouchterlony, O. et al., Chap. 19 in: *Handbook of Experimental Immunology*, Wier, D., ed, Blackwell (1973). Plateau concentration of antibody is usually in the range of 0.1 to 0.2 mg/ml of serum (about 12M). Affinity of the antisera for the antigen is determined by preparing competitive binding curves, as described, for example, by Fisher, D., Chap. 42 in: *Manual of Clinical Immunology*, second edition, Rose and Friedman, eds., Amer. Soc. For Microbiology, Washington, D. C. (1980)

30 Antibody preparations prepared according to either protocol are useful in quantitative immunoassays which determine concentrations of antigen-bearing substances in biological samples; they are also used semi- quantitatively or qualitatively to identify the presence of antigen in a biological sample. In addition, antibodies are useful in various animal models of pneumococcal disease as a means of evaluating the protein used to make the antibody as a

potential vaccine target or as a means of evaluating the antibody as a potential immunotherapeutic or immunoprophylactic reagent.

4. Preparation of PCR Primers and Amplification of DNA

5 Various fragments of the *T. pallidum* genome, such as those of Tables 1-3 and SEQ ID NOS: 1-744 can be used, in accordance with the present invention, to prepare PCR primers for a variety of uses. The PCR primers are preferably at least 15 bases, and more preferably at least 18 bases in length. When selecting a primer sequence, it is preferred that the primer pairs have approximately the same G/C ratio, so that melting temperatures are approximately the same. The
10 PCR primers and amplified DNA of this Example find use in the Examples that follow.

5. Isolation of a Selected DNA Clone From *T. pallidum*

Three approaches are used to isolate a *T. pallidum* clone comprising a polynucleotide of the present invention from any *T. pallidum* genomic DNA library. The *T. pallidum* strain
15 B31PU has been deposited as a convenient source for obtaining a *T. pallidum* strain although a wide variety of strains *T. pallidum* strains can be used which are known in the art.

10 *T. pallidum* genomic DNA is prepared using the following method. A 20ml overnight bacterial culture grown in a rich medium (e.g., Trypticase Soy Broth, Brain Heart Infusion broth or Super broth), pelleted, ished two times with TES (30mM Tris-pH 8.0, 25mM EDTA, 50mM NaCl), and resuspended in 5ml high salt TES (2.5M NaCl). Lysostaphin is added to final concentration of approx 50ug/ml and the mixture is rotated slowly 1 hour at 37C to make protoplast cells. The solution is then placed in incubator (or place in a shaking water bath) and warmed to 55C. Five hundred micro liter of 20% sarcosyl in TES (final concentration 2%) is then added to lyse the cells. Next, guanidine HCl is added to a final concentration of 7M (3.69g
20 in 5.5 ml). The mixture is swirled slowly at 55C for 60-90 min (solution should clear). A CsCl gradient is then set up in SW41 ultra clear tubes using 2.0ml 5.7M CsCl and overlaying with 2.85M CsCl. The gradient is carefully overlayed with the DNA-containing GuHCl solution. The gradient is spun at 30,000 rpm, 20C for 24 hr and the lower DNA band is collected. The volume is increased to 5 ml with TE buffer. The DNA is then treated with protease K (10 ug/ml)
25 overnight at 37 C, and precipitated with ethanol. The precipitated DNA is resuspended in a desired buffer.
30

In the first method, a plasmid is directly isolated by screening a plasmid *T. pallidum* genomic DNA library using a polynucleotide probe corresponding to a polynucleotide of the present invention. Particularly, a specific polynucleotide with 30-40 nucleotides is synthesized
35 using an Applied Biosystems DNA synthesizer according to the sequence reported. The oligonucleotide is labeled, for instance, with ^{32}P - γ -ATP using T4 polynucleotide kinase and purified according to routine methods. (See, e.g., Maniatis et al., Molecular Cloning: A Laboratory Manual, Cold Spring Harbor Press, Cold Spring, NY (1982).) The library is

transformed into a suitable host, as indicated above (such as XL-1 Blue (Stratagene)) using techniques known to those of skill in the art. *See, e.g.*, Sambrook et al. MOLECULAR CLONING: A LABORATORY MANUAL (Cold Spring Harbor, N.Y. 2nd ed. 1989); Ausubel et al., CURRENT PROTOCOLS IN MOLECULAR BIOLOGY (John Wiley and Sons, N.Y. 1989). The transformants are plated on 1.5% agar plates (containing the appropriate selection agent, e.g., ampicillin) to a density of about 150 transformants (colonies) per plate. These plates are screened using Nylon membranes according to routine methods for bacterial colony screening. *See, e.g.*, Sambrook et al. MOLECULAR CLONING: A LABORATORY MANUAL (Cold Spring Harbor, N.Y. 2nd ed. 1989); Ausubel et al., CURRENT PROTOCOLS IN MOLECULAR BIOLOGY (John Wiley and Sons, N.Y. 1989) or other techniques known to those of skill in the art.

Alternatively, two primers of 15-25 nucleotides derived from the 5' and 3' ends of a polynucleotide of SEQ ID NOS:1-744 are synthesized and used to amplify the desired DNA by PCR using a *T. pallidum* genomic DNA prep as a template. PCR is carried out under routine conditions, for instance, in 25 µl of reaction mixture with 0.5 ug of the above DNA template. A convenient reaction mixture is 1.5-5 mM MgCl₂, 0.01% (w/v) gelatin, 20 µM each of dATP, dCTP, dGTP, dTTP, 25 pmol of each primer and 0.25 Unit of Taq polymerase. Thirty five cycles of PCR (denaturation at 94°C for 1 min; annealing at 55°C for 1 min; elongation at 72°C for 1 min) are performed with a Perkin-Elmer Cetus automated thermal cycler. The amplified product is analyzed by agarose gel electrophoresis and the DNA band with expected molecular weight is excised and purified. The PCR product is verified to be the selected sequence by subcloning and sequencing the DNA product.

Finally, overlapping oligos of the DNA sequences of SEQ ID NOS:1-744 can be chemically synthesized and used to generate a nucleotide sequence of desired length using PCR methods known in the art.

6(a). Expression and Purification Borrelia polypeptides in E. coli

The bacterial expression vector pQE60 is used for bacterial expression of some of the polypeptide fragments of the present invention. (QIAGEN, Inc., 9259 Eton Avenue, Chatsworth, CA, 91311). pQE60 encodes ampicillin antibiotic resistance ("Ampr") and contains a bacterial origin of replication ("ori"), an IPTG inducible promoter, a ribosome binding site ("RBS"), six codons encoding histidine residues that allow affinity purification using nickel-nitrilo-tri-acetic acid ("Ni-NTA") affinity resin (QIAGEN, Inc., *supra*) and suitable single restriction enzyme cleavage sites. These elements are arranged such that an inserted DNA fragment encoding a polypeptide expresses that polypeptide with the six His residues (i.e., a "6 X His tag") covalently linked to the carboxyl terminus of that polypeptide.

The DNA sequence encoding the desired portion of a *T. pallidum* protein of the present invention is amplified from *T. pallidum* genomic DNA using PCR oligonucleotide primers

which anneal to the 5' and 3' sequences coding for the portions of the *T. pallidum* polynucleotide shown in SEQ ID NOS:1-744. Additional nucleotides containing restriction sites to facilitate cloning in the pQE60 vector are added to the 5' and 3' sequences, respectively.

For cloning the mature protein, the 5' primer has a sequence containing an appropriate restriction site followed by nucleotides of the amino terminal coding sequence of the desired *T. pallidum* polynucleotide sequence in SEQ ID NOS:1-744. One of ordinary skill in the art would appreciate that the point in the protein coding sequence where the 5' and 3' primers begin may be varied to amplify a DNA segment encoding any desired portion of the complete protein shorter or longer than the mature form. The 3' primer has a sequence containing an appropriate restriction site followed by nucleotides complementary to the 3' end of the polypeptide coding sequence of SEQ ID NOS:1-744, excluding a stop codon, with the coding sequence aligned with the restriction site so as to maintain its reading frame with that of the six His codons in the pQE60 vector.

The amplified *T. pallidum* DNA fragment and the vector pQE60 are digested with restriction enzymes which recognize the sites in the primers and the digested DNAs are then ligated together. The *T. pallidum* DNA is inserted into the restricted pQE60 vector in a manner which places the *T. pallidum* protein coding region downstream from the IPTG-inducible promoter and in-frame with an initiating AUG and the six histidine codons.

The ligation mixture is transformed into competent *E. coli* cells using standard procedures such as those described by Sambrook et al., *supra*. *E. coli* strain M15/rep4, containing multiple copies of the plasmid pREP4, which expresses the lac repressor and confers kanamycin resistance ("Kanr"), is used in carrying out the illustrative example described herein. This strain, which is only one of many that are suitable for expressing a *T. pallidum* polypeptide, is available commercially (QIAGEN, Inc., *supra*). Transformants are identified by their ability to grow on LB agar plates in the presence of ampicillin and kanamycin. Plasmid DNA is isolated from resistant colonies and the identity of the cloned DNA confirmed by restriction analysis, PCR and DNA sequencing.

Clones containing the desired constructs are grown overnight ("O/N") in liquid culture in LB media supplemented with both ampicillin (100 µg/ml) and kanamycin (25 µg/ml). The O/N culture is used to inoculate a large culture, at a dilution of approximately 1:25 to 1:250. The cells are grown to an optical density at 600 nm ("OD600") of between 0.4 and 0.6. Isopropyl-β-D-thiogalactopyranoside ("IPTG") is then added to a final concentration of 1 mM to induce transcription from the lac repressor sensitive promoter, by inactivating the lacI repressor. Cells subsequently are incubated further for 3 to 4 hours. Cells then are harvested by centrifugation.

The cells are then stirred for 3-4 hours at 4°C in 6M guanidine-HCl, pH 8. The cell debris is removed by centrifugation, and the supernatant containing the *T. pallidum* polypeptide is loaded onto a nickel-nitrilo-tri-acetic acid ("Ni-NTA") affinity resin column (QIAGEN, Inc., *supra*). Proteins with a 6 x His tag bind to the Ni-NTA resin with high affinity are purified in a

simple one-step procedure (for details see: The QIAexpressionist, 1995, QIAGEN, Inc., *supra*). Briefly the supernatant is loaded onto the column in 6 M guanidine-HCl, pH 8, the column is first washed with 10 volumes of 6 M guanidine-HCl, pH 8, then washed with 10 volumes of 6 M guanidine-HCl pH 6, and finally the *T. pallidum* polypeptide is eluted with 6 M guanidine-HCl, pH 5.

The purified protein is then renatured by dialyzing it against phosphate-buffered saline (PBS) or 50 mM Na-acetate, pH 6 buffer plus 200 mM NaCl. Alternatively, the protein could be successfully refolded while immobilized on the Ni-NTA column. The recommended conditions are as follows: renature using a linear 6M-1M urea gradient in 500 mM NaCl, 20% glycerol, 20 mM Tris/HCl pH 7.4, containing protease inhibitors. The renaturation should be performed over a period of 1.5 hours or more. After renaturation the proteins can be eluted by the addition of 250 mM immidazole. Immidazole is removed by a final dialyzing step against PBS or 50 mM sodium acetate pH 6 buffer plus 200 mM NaCl. The purified protein is stored at 4°C or frozen at -80°C.

The polypeptide of the present invention are also prepared using a non-denaturing protein purification method. For these polypeptides, the cell pellet from each liter of culture is resuspended in 25 mls of Lysis Buffer A at 4°C (Lysis Buffer A = 50 mM Na-phosphate, 300 mM NaCl, 10 mM 2-mercaptoethanol, 10% Glycerol, pH 7.5 with 1 tablet of Complete EDTA-free protease inhibitor cocktail (Boehringer Mannheim #1873580) per 50 ml of buffer). Absorbance at 550 nm is approximately 10-20 O.D./ml. The suspension is then put through three freeze/thaw cycles from -70°C (using a ethanol-dry ice bath) up to room temperature. The cells are lysed via sonication in short 10 sec bursts over 3 minutes at approximately 80W while kept on ice. The sonicated sample is then centrifuged at 15,000 RPM for 30 minutes at 4°C. The supernatant is passed through a column containing 1.0 ml of CL-4B resin to pre-clear the sample of any proteins that may bind to agarose non-specifically, and the flow-through fraction is collected.

The pre-cleared flow-through is applied to a nickel-nitrilo-tri-acetic acid ("Ni-NTA") affinity resin column (Quiagen, Inc., *supra*). Proteins with a 6 X His tag bind to the Ni-NTA resin with high affinity and can be purified in a simple one-step procedure. Briefly, the supernatant is loaded onto the column in Lysis Buffer A at 4°C, the column is first washed with 10 volumes of Lysis Buffer A until the A280 of the eluate returns to the baseline. Then, the column is washed with 5 volumes of 40 mM Imidazole (92% Lysis Buffer A / 8% Buffer B) (Buffer B = 50 mM Na-Phosphate, 300 mM NaCl, 10% Glycerol, 10 mM 2-mercaptoethanol, 500 mM Imidazole, pH of the final buffer should be 7.5). The protein is eluted off of the column with a series of increasing Imidazole solutions made by adjusting the ratios of Lysis Buffer A to Buffer B. Three different concentrations are used: 3 volumes of 75 mM Imidazole, 3 volumes of 150 mM Imidazole, 5 volumes of 500 mM Imidazole. The fractions containing the purified protein are analyzed using 8 %, 10 % or 14% SDS-PAGE depending on the protein size. The purified protein is then dialyzed 2X against phosphate-buffered saline (PBS) in order to place it

into an easily workable buffer. The purified protein is stored at 4°C or frozen at -80°.

The following alternative method may be used to purify *T. pallidum* expressed in *E. coli* when it is present in the form of inclusion bodies. Unless otherwise specified, all of the following steps are conducted at 4-10°C.

5 Upon completion of the production phase of the *E. coli* fermentation, the cell culture is cooled to 4-10°C and the cells are harvested by continuous centrifugation at 15,000 rpm (Heraeus Sepatech). On the basis of the expected yield of protein per unit weight of cell paste and the amount of purified protein required, an appropriate amount of cell paste, by weight, is suspended in a buffer solution containing 100 mM Tris, 50 mM EDTA, pH 7.4. The cells are
10 dispersed to a homogeneous suspension using a high shear mixer.

15 The cells are then lysed by passing the solution through a microfluidizer (Microfluidics, Corp. or APV Gaulin, Inc.) twice at 4000-6000 psi. The homogenate is then mixed with NaCl solution to a final concentration of 0.5 M NaCl, followed by centrifugation at 7000 x g for 15 min. The resultant pellet is washed again using 0.5M NaCl, 100 mM Tris, 50 mM EDTA; pH
7.4.

15 The resulting washed inclusion bodies are solubilized with 1.5 M guanidine hydrochloride (GuHCl) for 2-4 hours. After 7000 x g centrifugation for 15 min., the pellet is discarded and the *T. pallidum* polypeptide-containing supernatant is incubated at 4°C overnight to allow further GuHCl extraction.

20 Following high speed centrifugation (30,000 x g) to remove insoluble particles, the GuHCl solubilized protein is refolded by quickly mixing the GuHCl extract with 20 volumes of buffer containing 50 mM sodium, pH 4.5, 150 mM NaCl, 2 mM EDTA by vigorous stirring. The refolded diluted protein solution is kept at 4°C without mixing for 12 hours prior to further purification steps.

25 To clarify the refolded *T. pallidum* polypeptide solution, a previously prepared tangential filtration unit equipped with 0.16 µm membrane filter with appropriate surface area (e.g., Filtron), equilibrated with 40 mM sodium acetate, pH 6.0 is employed. The filtered sample is loaded onto a cation exchange resin (e.g., Poros HS-50, Perseptive Biosystems). The column is washed with 40 mM sodium acetate, pH 6.0 and eluted with 250 mM, 500 mM, 1000 mM, and
30 1500 mM NaCl in the same buffer, in a stepwise manner. The absorbance at 280 nm of the effluent is continuously monitored. Fractions are collected and further analyzed by SDS-PAGE.

35 Fractions containing the *T. pallidum* polypeptide are then pooled and mixed with 4 volumes of water. The diluted sample is then loaded onto a previously prepared set of tandem columns of strong anion (Poros HQ-50, Perseptive Biosystems) and weak anion (Poros CM-20, Perseptive Biosystems) exchange resins. The columns are equilibrated with 40 mM sodium acetate, pH 6.0. Both columns are washed with 40 mM sodium acetate, pH 6.0, 200 mM NaCl.

The CM-20 column is then eluted using a 10 column volume linear gradient ranging from 0.2 M NaCl, 50 mM sodium acetate, pH 6.0 to 1.0 M NaCl, 50 mM sodium acetate, pH 6.5. Fractions are collected under constant A_{280} monitoring of the effluent. Fractions containing the *T. pallidum* polypeptide (determined, for instance, by 16% SDS-PAGE) are then pooled.

5 The resultant *T. pallidum* polypeptide exhibits greater than 95% purity after the above refolding and purification steps. No major contaminant bands are observed from Commassie blue stained 16% SDS-PAGE gel when 5 μ g of purified protein is loaded. The purified protein is also tested for endotoxin/LPS contamination, and typically the LPS content is less than 0.1 ng/ml according to LAL assays.

10 **6(b). Alternative Expression and Purification Borrelia polypeptides in E. coli**

15 The vector pQE10 is alternatively used to clone and express some of the polypeptides of the present invention for use in the soft tissue and systemic infection models discussed below.

15 The difference being such that an inserted DNA fragment encoding a polypeptide expresses that polypeptide with the six His residues (i.e., a "6 X His tag") covalently linked to the amino terminus of that polypeptide. The bacterial expression vector pQE10 (QIAGEN, Inc., 9259 Eton Avenue, Chatsworth, CA, 91311) was used in this example. The components of the pQE10 plasmid are arranged such that the inserted DNA sequence encoding a polypeptide of the present

20 invention expresses the polypeptide with the six His residues (i.e., a "6 X His tag") covalently linked to the amino terminus.

25 The DNA sequences encoding the desired portions of a polypeptide of SEQ ID NOS:1-744 were amplified using PCR oligonucleotide primers from genomic *T. pallidum* DNA. The PCR primers anneal to the nucleotide sequences encoding the desired amino acid sequence of a polypeptide of the present invention. Additional nucleotides containing restriction sites to facilitate cloning in the pQE10 vector were added to the 5' and 3' primer sequences, respectively.

30 For cloning a polypeptide of the present invention, the 5' and 3' primers were selected to amplify their respective nucleotide coding sequences. One of ordinary skill in the art would appreciate that the point in the protein coding sequence where the 5' and 3' primers begins may be varied to amplify a DNA segment encoding any desired portion of a polypeptide of the present invention. The 5' primer was designed so the coding sequence of the 6 X His tag is aligned with the restriction site so as to maintain its reading frame with that of *T. pallidum* polypeptide. The 3' was designed to include an stop codon. The amplified DNA fragment was then cloned, and the protein expressed, as described above for the pQE60 plasmid.

35 The DNA sequences encoding the amino acid sequences of SEQ ID NOS:1-744 may also be cloned and expressed as fusion proteins by a protocol similar to that described directly above, wherein the pET-32b(+) vector (Novagen, 601 Science Drive, Madison, WI 53711) is preferentially used in place of pQE10.

The above methods are not limited to the polypeptide fragments actually produced. The above method, like the methods below, can be used to produce either full length polypeptides or desired fragments thereof.

5 **6(c). Alternative Expression and Purification of Borrelia polypeptides in
E. coli**

The bacterial expression vector pQE60 is used for bacterial expression in this example (QIAGEN, Inc., 9259 Eton Avenue, Chatsworth, CA, 91311). However, in this example, the polypeptide coding sequence is inserted such that translation of the six His codons is prevented 10 and, therefore, the polypeptide is produced with no 6 X His tag.

The DNA sequence encoding the desired portion of the *T. pallidum* amino acid sequence is amplified from an *T. pallidum* genomic DNA prep the deposited DNA clones using PCR oligonucleotide primers which anneal to the 5' and 3' nucleotide sequences corresponding to the desired portion of the *T. pallidum* polypeptides. Additional nucleotides containing restriction 15 sites to facilitate cloning in the pQE60 vector are added to the 5' and 3' primer sequences.

For cloning a *T. pallidum* polypeptides of the present invention, 5' and 3' primers are selected to amplify their respective nucleotide coding sequences. One of ordinary skill in the art would appreciate that the point in the protein coding sequence where the 5' and 3' primers begin may be varied to amplify a DNA segment encoding any desired portion of a polypeptide of the 20 present invention. The 3' and 5' primers contain appropriate restriction sites followed by nucleotides complementary to the 5' and 3' ends of the coding sequence respectively. The 3' primer is additionally designed to include an in-frame stop codon.

The amplified *T. pallidum* DNA fragments and the vector pQE60 are digested with restriction enzymes recognizing the sites in the primers and the digested DNAs are then ligated 25 together. Insertion of the *T. pallidum* DNA into the restricted pQE60 vector places the *T. pallidum* protein coding region including its associated stop codon downstream from the IPTG-inducible promoter and in-frame with an initiating AUG. The associated stop codon prevents translation of the six histidine codons downstream of the insertion point.

The ligation mixture is transformed into competent *E. coli* cells using standard procedures 30 such as those described by Sambrook et al. *E. coli* strain M15/rep4, containing multiple copies of the plasmid pREP4, which expresses the lac repressor and confers kanamycin resistance ("Kanr"), is used in carrying out the illustrative example described herein. This strain, which is only one of many that are suitable for expressing *T. pallidum* polypeptide, is available commercially (QIAGEN, Inc., *supra*). Transformants are identified by their ability to grow on 35 LB plates in the presence of ampicillin and kanamycin. Plasmid DNA is isolated from resistant colonies and the identity of the cloned DNA confirmed by restriction analysis, PCR and DNA sequencing.

Clones containing the desired constructs are grown overnight ("O/N") in liquid culture in LB media supplemented with both ampicillin (100 µg/ml) and kanamycin (25 µg/ml). The O/N

culture is used to inoculate a large culture, at a dilution of approximately 1:25 to 1:250. The cells are grown to an optical density at 600 nm ("OD600") of between 0.4 and 0.6. Isopropyl-β-D-thiogalactopyranoside ("IPTG") is then added to a final concentration of 1 mM to induce transcription from the *lac* repressor sensitive promoter, by inactivating the lacI repressor. Cells subsequently are incubated further for 3 to 4 hours. Cells then are harvested by centrifugation.

To purify the *T. pallidum* polypeptide, the cells are then stirred for 3-4 hours at 4°C in 6M guanidine-HCl, pH 8. The cell debris is removed by centrifugation, and the supernatant containing the *T. pallidum* polypeptide is dialyzed against 50 mM Na-acetate buffer pH 6, supplemented with 200 mM NaCl. Alternatively, the protein can be successfully refolded by dialyzing it against 500 mM NaCl, 20% glycerol, 25 mM Tris/HCl pH 7.4, containing protease inhibitors. After renaturation the protein can be purified by ion exchange, hydrophobic interaction and size exclusion chromatography. Alternatively, an affinity chromatography step such as an antibody column can be used to obtain pure *T. pallidum* polypeptide. The purified protein is stored at 4°C or frozen at -80°C.

The following alternative method may be used to purify *T. pallidum* polypeptides expressed in *E. coli* when it is present in the form of inclusion bodies. Unless otherwise specified, all of the following steps are conducted at 4-10°C.

Upon completion of the production phase of the *E. coli* fermentation, the cell culture is cooled to 4-10°C and the cells are harvested by continuous centrifugation at 15,000 rpm (Heraeus Sepatech). On the basis of the expected yield of protein per unit weight of cell paste and the amount of purified protein required, an appropriate amount of cell paste, by weight, is suspended in a buffer solution containing 100 mM Tris, 50 mM EDTA, pH 7.4. The cells are dispersed to a homogeneous suspension using a high shear mixer.

The cells were then lysed by passing the solution through a microfluidizer (Microfluidics, Corp. or APV Gaulin, Inc.) twice at 4000-6000 psi. The homogenate is then mixed with NaCl solution to a final concentration of 0.5 M NaCl, followed by centrifugation at 7000 x g for 15 min. The resultant pellet is washed again using 0.5M NaCl, 100 mM Tris, 50 mM EDTA, pH 7.4.

The resulting washed inclusion bodies are solubilized with 1.5 M guanidine hydrochloride (GuHCl) for 2-4 hours. After 7000 x g centrifugation for 15 min., the pellet is discarded and the *T. pallidum* polypeptide-containing supernatant is incubated at 4°C overnight to allow further GuHCl extraction.

Following high speed centrifugation (30,000 x g) to remove insoluble particles, the GuHCl solubilized protein is refolded by quickly mixing the GuHCl extract with 20 volumes of buffer containing 50 mM sodium, pH 4.5, 150 mM NaCl, 2 mM EDTA by vigorous stirring.

The refolded diluted protein solution is kept at 4°C without mixing for 12 hours prior to further

purification steps.

To clarify the refolded *T. pallidum* polypeptide solution, a previously prepared tangential filtration unit equipped with 0.16 µm membrane filter with appropriate surface area (e.g., Filtron), equilibrated with 40 mM sodium acetate, pH 6.0 is employed. The filtered sample is loaded onto a cation exchange resin (e.g., Poros HS-50, Perseptive Biosystems). The column is washed with 40 mM sodium acetate, pH 6.0 and eluted with 250 mM, 500 mM, 1000 mM, and 1500 mM NaCl in the same buffer, in a stepwise manner. The absorbance at 280 nm of the effluent is continuously monitored. Fractions are collected and further analyzed by SDS-PAGE.

Fractions containing the *T. pallidum* polypeptide are then pooled and mixed with 4 volumes of water. The diluted sample is then loaded onto a previously prepared set of tandem columns of strong anion (Poros HQ-50, Perseptive Biosystems) and weak anion (Poros CM-20, Perseptive Biosystems) exchange resins. The columns are equilibrated with 40 mM sodium acetate, pH 6.0. Both columns are washed with 40 mM sodium acetate, pH 6.0, 200 mM NaCl. The CM-20 column is then eluted using a 10 column volume linear gradient ranging from 0.2 M NaCl, 50 mM sodium acetate, pH 6.0 to 1.0 M NaCl, 50 mM sodium acetate, pH 6.5. Fractions are collected under constant A₂₈₀ monitoring of the effluent. Fractions containing the *T. pallidum* polypeptide (determined, for instance, by 16% SDS-PAGE) are then pooled.

The resultant *T. pallidum* polypeptide exhibits greater than 95% purity after the above refolding and purification steps. No major contaminant bands are observed from Commassie blue stained 16% SDS-PAGE gel when 5 µg of purified protein is loaded. The purified protein is also tested for endotoxin/LPS contamination, and typically the LPS content is less than 0.1 ng/ml according to LAL assays.

6(d). Cloning and Expression of *T. pallidum* in Other Bacteria

T. pallidum polypeptides can also be produced in: *T. pallidum* using the methods of S. Skinner et al., (1988) Mol. Microbiol. 2:289-297 or J. I. Moreno (1996) Protein Expr. Purif. 8(3):332-340; *Lactobacillus* using the methods of C. Rush et al., 1997 Appl. Microbiol. Biotechnol. 47(5):537-542; or in *Bacillus subtilis* using the methods Chang et al., U.S. Patent No. 4,952,508.

30 7. Cloning and Expression in COS Cells

A *T. pallidum* expression plasmid is made by cloning a portion of the DNA encoding a *T. pallidum* polypeptide into the expression vector pDNAI/Amp or pDNAII (which can be obtained from Invitrogen, Inc.). The expression vector pDNAI/amp contains: (1) an *E. coli* origin of replication effective for propagation in *E. coli* and other prokaryotic cells; (2) an ampicillin resistance gene for selection of plasmid-containing prokaryotic cells; (3) an SV40 origin of replication for propagation in eukaryotic cells; (4) a CMV promoter, a polylinker, an SV40 intron; (5) several codons encoding a hemagglutinin fragment (i.e., an "HA" tag to

facilitate purification) followed by a termination codon and polyadenylation signal arranged so that a DNA can be conveniently placed under expression control of the CMV promoter and operably linked to the SV40 intron and the polyadenylation signal by means of restriction sites in the polylinker. The HA tag corresponds to an epitope derived from the influenza hemagglutinin 5 protein described by Wilson et al. 1984 Cell 37:767. The fusion of the HA tag to the target protein allows easy detection and recovery of the recombinant protein with an antibody that recognizes the HA epitope. pDNAIII contains, in addition, the selectable neomycin marker.

A DNA fragment encoding a *T. pallidum* polypeptide is cloned into the polylinker region of the vector so that recombinant protein expression is directed by the CMV promoter. The 10 plasmid construction strategy is as follows. The DNA from a *T. pallidum* genomic DNA prep is amplified using primers that contain convenient restriction sites, much as described above for construction of vectors for expression of *T. pallidum* in *E. coli*. The 5' primer contains a Kozak sequence, an AUG start codon, and nucleotides of the 5' coding region of the *T. pallidum* polypeptide. The 3' primer, contains nucleotides complementary to the 3' coding sequence of the 15 *T. pallidum* DNA, a stop codon, and a convenient restriction site.

The PCR amplified DNA fragment and the vector, pDNAI/Amp, are digested with appropriate restriction enzymes and then ligated. The ligation mixture is transformed into an appropriate *E. coli* strain such as SURE™ (Stratagene Cloning Systems, La Jolla, CA 92037), and the transformed culture is plated on ampicillin media plates which then are incubated to allow 20 growth of ampicillin resistant colonies. Plasmid DNA is isolated from resistant colonies and examined by restriction analysis or other means for the presence of the fragment encoding the *T. pallidum* polypeptide.

For expression of a recombinant *T. pallidum* polypeptide, COS cells are transfected with an expression vector, as described above, using DEAE-dextran, as described, for instance, by 25 Sambrook et al. (*supra*). Cells are incubated under conditions for expression of *T. pallidum* by the vector.

Expression of the *T. pallidum*-HA fusion protein is detected by radiolabeling and immunoprecipitation, using methods described in, for example Harlow et al., *supra*. To this end, two days after transfection, the cells are labeled by incubation in media containing ³⁵S-30 cysteine for 8 hours. The cells and the media are collected, and the cells are washed and the lysed with detergent-containing RIPA buffer: 150 mM NaCl, 1% NP-40, 0.1% SDS, 1% NP-40, 0.5% DOC, 50 mM TRIS, pH 7.5, as described by Wilson et al. (*supra*). Proteins are precipitated from the cell lysate and from the culture media using an HA-specific monoclonal antibody. The precipitated proteins then are analyzed by SDS-PAGE and autoradiography. An 35 expression product of the expected size is seen in the cell lysate, which is not seen in negative controls.

8. Cloning and Expression in CHO Cells

The vector pC4 is used for the expression of *T. pallidum* polypeptide in this example. Plasmid pC4 is a derivative of the plasmid pSV2-dhfr (ATCC Accession No. 37146). The plasmid contains the mouse DHFR gene under control of the SV40 early promoter. Chinese hamster ovary cells or other cells lacking dihydrofolate activity that are transfected with these 5 plasmids can be selected by growing the cells in a selective medium (alpha minus MEM, Life Technologies) supplemented with the chemotherapeutic agent methotrexate. The amplification of the DHFR genes in cells resistant to methotrexate (MTX) has been well documented. See, e.g., Alt et al., 1978, J. Biol. Chem. 253:1357-1370; Hamlin et al., 1990, Biochem. et Biophys. Acta, 1097:107-143; Page et al., 1991, Biotechnology 9:64-68. Cells grown in increasing 10 concentrations of MTX develop resistance to the drug by overproducing the target enzyme, DHFR, as a result of amplification of the DHFR gene. If a second gene is linked to the DHFR gene, it is usually co-amplified and over-expressed. It is known in the art that this approach may be used to develop cell lines carrying more than 1,000 copies of the amplified gene(s). Subsequently, when the methotrexate is withdrawn, cell lines are obtained which contain the 15 amplified gene integrated into one or more chromosome(s) of the host cell.

Plasmid pC4 contains the strong promoter of the long terminal repeat (LTR) of the Rous Sarcoma Virus, for expressing a polypeptide of interest, Cullen, et al. (1985) Mol. Cell. Biol. 5:438-447; plus a fragment isolated from the enhancer of the immediate early gene of human cytomegalovirus (CMV), Boshart, et al., 1985, Cell 41:521-530. Downstream of the promoter 20 are the following single restriction enzyme cleavage sites that allow the integration of the genes: *Bam* HI, *Xba* I, and *Asp* 718. Behind these cloning sites the plasmid contains the 3' intron and polyadenylation site of the rat preproinsulin gene. Other high efficiency promoters can also be used for the expression, e.g., the human β -actin promoter, the SV40 early or late promoters or the long terminal repeats from other retroviruses, e.g., HIV and HTLV. Clontech's Tet-Off and 25 Tet-On gene expression systems and similar systems can be used to express the *T. pallidum* polypeptide in a regulated way in mammalian cells (Gossen et al., 1992, Proc. Natl. Acad. Sci. USA 89:5547-5551. For the polyadenylation of the mRNA other signals, e.g., from the human growth hormone or globin genes can be used as well. Stable cell lines carrying a gene of interest integrated into the chromosomes can also be selected upon co-transfection with a selectable 30 marker such as gpt, G418 or hygromycin. It is advantageous to use more than one selectable marker in the beginning, e.g., G418 plus methotrexate.

The plasmid pC4 is digested with the restriction enzymes and then dephosphorylated using calf intestinal phosphates by procedures known in the art. The vector is then isolated from a 1% agarose gel. The DNA sequence encoding the *T. pallidum* polypeptide is amplified using 35 PCR oligonucleotide primers corresponding to the 5' and 3' sequences of the desired portion of the gene. A 5' primer containing a restriction site, a Kozak sequence, an AUG start codon, and nucleotides of the 5' coding region of the *T. pallidum* polypeptide is synthesized and used. A 3' primer, containing a restriction site, stop codon, and nucleotides complementary to the 3' coding sequence of the *T. pallidum* polypeptides is synthesized and used. The amplified fragment is

digested with the restriction endonucleases and then purified again on a 1% agarose gel. The isolated fragment and the dephosphorylated vector are then ligated with T4 DNA ligase. *E. coli* HB101 or XL-1 Blue cells are then transformed and bacteria are identified that contain the fragment inserted into plasmid pC4 using, for instance, restriction enzyme analysis.

- 5 Chinese hamster ovary cells lacking an active DHFR gene are used for transfection. Five µg of the expression plasmid pC4 is cotransfected with 0.5 µg of the plasmid pSVneo using a lipid-mediated transfection agent such as Lipofectin™ or LipofectAMINE™ (LifeTechnologies Gaithersburg, MD). The plasmid pSV2-neo contains a dominant selectable marker, the *neo* gene from Tn5 encoding an enzyme that confers resistance to a group of antibiotics including G418.
- 10 The cells are seeded in alpha minus MEM supplemented with 1 mg/ml G418. After 2 days, the cells are trypsinized and seeded in hybridoma cloning plates (Greiner, Germany) in alpha minus MEM supplemented with 10, 25, or 50 ng/ml of methotrexate plus 1 mg/ml G418. After about 10-14 days single clones are trypsinized and then seeded in 6-well petri dishes or 10 ml flasks using different concentrations of methotrexate (50 nM, 100 nM, 200 nM, 400 nM, 800 nM).
- 15 Clones growing at the highest concentrations of methotrexate are then transferred to new 6-well plates containing even higher concentrations of methotrexate (1 µM, 2 µM, 5 µM, 10 mM, 20 mM). The same procedure is repeated until clones are obtained which grow at a concentration of 100-200 µM. Expression of the desired gene product is analyzed, for instance, by SDS-PAGE and Western blot or by reversed phase HPLC analysis.
- 20 The disclosure of all publications (including patents, patent applications, journal articles, laboratory manuals, books, or other documents) cited herein are hereby incorporated by reference in their entireties SEQ ID NOS: 1-744 are hereby incorporated into the specification by reference.
- 25 The present invention is not to be limited in scope by the specific embodiments described herein, which are intended as single illustrations of individual aspects of the invention.
- Functionally equivalent methods and components are within the scope of the invention, in addition to those shown and described herein and will become apparent to those skilled in the art from the foregoing description and accompanying drawings. Such modifications are intended to fall within the scope of the appended claims.

TABLE 1.

Treponema pallidum - Coding regions containing known sequences

Contig ID	Orf ID	Start (nt)	Stop (nt)	match accession	match gene name	percent ident	HSP nt length
1	1	219	4	gbIU55214I	Treponema pallidum RuvB' gene, partial cds, and TroA, TroB, TroC, TroD, TroR, Pgm, Pf, and Tpp15 genes, complete cds	93	202
1	2	110	493	gbIU55214I	Treponema pallidum RuvB' gene, partial cds, and TroA, TroB, TroC, TroD, TroR, Pgm, Pf, and Tpp15 genes, complete cds	97	199
1	3	1167	226	gbIU55214I	Treponema pallidum RuvB' gene, partial cds, and TroA, TroB, TroC, TroD, TroR, Pgm, Pf, and Tpp15 genes, complete cds	99	829
1	4	1237	1644	gbIU55214I	Treponema pallidum RuvB' gene, partial cds, and TroA, TroB, TroC, TroD, TroR, Pgm, Pf, and Tpp15 genes, complete cds	98	285
3	1	225	2105	gbIU32683I	Treponema pallidum cytoplasmic filament protein A (cfpA) gene, complete cds	99	1800
16	44	25607	26131	embIX61228ITP33G	Treponema pallidum Tp33 gene (partial)	96	84
18	1	647	1471	gbIU36839I	Treponema pallidum putative switch protein FliM (fliM) gene, partial cds, and FliY (fliY) gene, complete cds, export apparatus proteins FliP (fliP), FliQ (fliQ), FliR (fliR) and FlihB (flihB), signal transducing receptor FlihA (flihA), GTP binding protein >	100	797
18	2	1687	2574	gbIU36839I	Treponema pallidum putative switch protein FliM (fliM) gene, partial cds, and FliY (fliY) gene, complete cds, export apparatus proteins FliP (fliP), FliQ (fliQ), FliR (fliR) and FlihB (flihB), signal transducing receptor FlihA (flihA), GTP binding protein >	98	843
18	3	3200	2589	gbIU36839I	Treponema pallidum putative switch protein FliM (fliM) gene, partial cds, and FliY (fliY) gene, complete cds, export apparatus proteins FliP (fliP),	100	459

TABLE 1.

Treponema pallidum - Coding regions containing known sequences

			FliQ (fliQ), FliR (fliR) and FlihB (flihB), signal transducing receptor FlihA (flihA), GTP binding protein >	
21	4	3095	2532 gbm3240II T.pallidum pallidum antigen TyFI gene, complete cds	99 426
21	5	3151	3807 gbm3240II T.pallidum pallidum antigen TyFI gene, complete cds	98 657
25	5	3042	2275 gbuU97563I Treponema pallidum FlaA homolog-1 and FlaA homolog-2 genes, complete cds	98 269
26	15	8262	9125 gbuU88957I Treponema pallidum major outer sheath protein homolog Msp (msp) gene, complete cds	68 417
26	16	10214	8373 gbuU88957I Treponema pallidum major outer sheath protein homolog Msp (msp) gene, complete cds	68 417
26	17	9123	9719 gbuU88957I Treponema pallidum major outer sheath protein homolog Msp (msp) gene, complete cds	85 122
26	19	11173	10226 gbuU88957I Treponema pallidum major outer sheath protein homolog Msp (msp) gene, complete cds	99 916
26	20	11397	11095 gbuU88957I Treponema pallidum major outer sheath protein homolog Msp (msp) gene, complete cds	97 228
26	21	12800	11316 gbuU88957I Treponema pallidum major outer sheath protein homolog Msp (msp) gene, complete cds	99 1079
26	22	11568	11777 gbuU88957I Treponema pallidum major outer sheath protein homolog Msp (msp) gene, complete cds	100 210
26	23	12824	13621 embIX57836I T.pallidum tmboC gene for membrane lipoprotein TPTMBCL	99 662
29	1	3	806 gbuU88957I Treponema pallidum major outer sheath protein homolog Msp (msp) gene, complete cds	64 407
29	2	884	165 gbuU88957I Treponema pallidum major outer sheath protein homolog Msp (msp) gene, complete cds	64 407
31	1	1	510 gbuU70661I Treponema pallidum GTP-binding protein, mccF-like protein and ATP-dependent DNA helicase (RecG) genes, complete cds	100 102

TABLE 1.

Treponema pallidum - Coding regions containing known sequences

32	1	3	458	gb M26525	T pallidum endoflagellar sheath protein (fiaA) gene, 3' end	100	423
32	2	358	810	gb U26525	T pallidum endoflagellar sheath protein (fiaA) gene, 3' end	100	384
34	7	5817	5134	gb U32683	Treponema pallidum cytoplasmic filament protein A (cfpA) gene, complete cds	98	65
35	20	12350	11691	gb U61534	Treponema pallidum 2-phospho-D-glycerate hydrolase (Eno) gene, complete cds	100	610
35	21	12918	12304	gb U61534	Treponema pallidum 2-phospho-D-glycerate hydrolase (Eno) gene, complete cds	99	604
35	22	13933	13091	gb U93844	Treponema pallidum lipoprotein homolog (tpN32) gene, complete cds, and 2-phospho-D-glycerate hydrolase (eno) gene, partial cds	100	780
35	23	14037	13792	gb U93844	Treponema pallidum lipoprotein homolog (tpN32) gene, complete cds, and 2-phospho-D-glycerate hydrolase (eno) gene, partial cds	99	169
35	24	14634	14050	gb U93844	Treponema pallidum lipoprotein homolog (tpN32) gene, complete cds, and 2-phospho-D-glycerate hydrolase (eno) gene, partial cds	100	544
35	25	15748	14915	gb U93844	Treponema pallidum lipoprotein homolog (tpN32) gene, complete cds, and 2-phospho-D-glycerate hydrolase (eno) gene, partial cds	100	834
35	26	15679	16056	gb U97358	Treponema pallidum 29K protein gene, complete cds	99	301
35	27	16719	15802	gb U97358	Treponema pallidum 29K protein gene, complete cds	99	178
36	1	526	1137	gb U57756	Treponema pallidum alanine racemase gene, partial cds	100	60
36	2	1157	22999	gb U57756	Treponema pallidum alanine racemase gene, partial cds	100	605
36	24	19592	20290	gb U97363	Treponema pallidum FiaA homolog-1 and FiaA homolog-2 genes, complete cds	98	261

TABLE 1.

Treponema pallidum - Coding regions containing known sequences

36	26	20754	20308	gbIU973631	Treponema pallidum FlmA homolog-1 and FlmA homolog-2 genes, complete cds	96	395
45	12	9409	10152	gbIM738251	Treponema pallidum 1-pyrroline-5-carboxylate reductase gene, complete cds	100	291
46	1	547	101	gbIU973611	Treponema pallidum H-ATPase homolog gene, partial cds	99	403
51	6	6397	7371	embIX612261	T.pallidum Tp75 gene (partial)	100	141
54	1	652	2	gbIU282191	Treponema pallidum flagellar hook (flgE), (orf4), flagellar motor protein (motA), flagellar motor protein (motB), (fliL), and flagellar switch protein (fliM) genes, complete cds, and (fliY) gene, partial cds	99	565
54	2	1306	905	gbIU282191	Treponema pallidum flagellar hook (flgE), (orf4), flagellar motor protein (motA), flagellar motor protein (motB), (fliL), and flagellar switch protein (fliM) genes, complete cds, and (fliY) gene, partial cds	99	388
54	3	1848	1237	gbIU282191	Treponema pallidum flagellar hook (flgE), (orf4), flagellar motor protein (motA), flagellar motor protein (motB), (fliL), and flagellar switch protein (fliM) genes, complete cds, and (fliY) gene, partial cds	100	503
54	4	1243	1479	gbIU282191	Treponema pallidum flagellar hook (flgE), (orf4), flagellar motor protein (motA), flagellar motor protein (motB), (fliL), and flagellar switch protein (fliM) genes, complete cds, and (fliY) gene, partial cds	100	194
54	5	2126	1755	gbIU282191	Treponema pallidum flagellar hook (flgE), (orf4), flagellar motor protein (motA), flagellar motor protein (motB), (fliL), and flagellar switch protein (fliM) genes, complete cds, and (fliY) gene, partial cds	96	192

TABLE 1.

Treponema pallidum - Coding regions containing known sequences

54	6	2329	1913	gbIU282191	Treponema pallidum flagellar hook (flgE), (orf4), flagellar motor protein (motA), flagellar motor protein (motB), (fllL), and flagellar switch protein (fliM) genes, complete cds, and (fliY) gene, partial cds	98	380
54	7	2735	2364	gbIU420121	Treponema pallidum putative treponemal aqueous protein Tap1 (tap1) gene, and flagellar hook assembly scaffolding protein (flgD) gene, complete cds	100	283
54	8	2872	2573	gbIU420121	Treponema pallidum putative treponemal aqueous protein Tap1 (tap1) gene, and flagellar hook assembly scaffolding protein (flgD) gene, complete cds	100	226
54	9	3388	2900	gbIU420121	Treponema pallidum putative treponemal aqueous protein Tap1 (tap1) gene, and flagellar hook assembly scaffolding protein (flgD) gene, complete cds	99	414
54	10	3169	4188	gbIU420121	Treponema pallidum putative treponemal aqueous protein Tap1 (tap1) gene, and flagellar hook assembly scaffolding protein (flgD) gene, complete cds	99	448
54	11	4191	3238	gbIU420121	Treponema pallidum putative treponemal aqueous protein Tap1 (tap1) gene, and flagellar hook assembly scaffolding protein (flgD) gene, complete cds	99	448
54	12	3683	3309	gbIU420121	Treponema pallidum putative treponemal aqueous protein Tap1 (tap1) gene, and flagellar hook assembly scaffolding protein (flgD) gene, complete cds	99	140
54	13	4529	4014	gbIU420121	Treponema pallidum putative treponemal aqueous protein Tap1 (tap1) gene, and flagellar hook assembly scaffolding protein (flgD) gene, complete cds	99	354

TABLE 1.

Treponema pallidum - Coding regions containing known sequences

54	14	4872	43211 gblU420121	Treponema pallidum putative treponemal aqueous protein Tap I (tap1) gene, and flagellar hook assembly scaffolding protein (flgD) gene, complete cds	99	122
54	15	4456	5550 gblU420121	Treponema pallidum putative treponemal aqueous protein Tap I (tap1) gene, and flagellar hook assembly scaffolding protein (flgD) gene, complete cds	99	122
55	3	759	1289 gblU973591	Treponema pallidum 76K protein gene, complete cds	100	96
55	4	1256	1753 gblU973591	Treponema pallidum 76K protein gene, complete cds	99	494
55	5	2336	3370 gblU973591	Treponema pallidum 76K protein gene, complete cds	99	919
56	1	1	225 gblU618511	Treponema pallidum histidine kinase CheA (cheA), and chemotaxis proteins CheW (cheW), CheX (cheX) and CheY (cheY) genes, complete cds	100	138
56	2	129	308 gblU618511	Treponema pallidum histidine kinase CheA (cheA), and chemotaxis proteins CheW (cheW), CheX (cheX) and CheY (cheY) genes, complete cds	100	167
56	3	281	1669 gblU618511	Treponema pallidum histidine kinase CheA (cheA), and chemotaxis proteins CheW (cheW), CheX (cheX) and CheY (cheY) genes, complete cds	100	1389
56	4	1667	2155 gblU618511	Treponema pallidum histidine kinase CheA (cheA), and chemotaxis proteins CheW (cheW), CheX (cheX) and CheY (cheY) genes, complete cds	100	489
56	5	2128	2601 gblU618511	Treponema pallidum histidine kinase CheA (cheA), and chemotaxis proteins CheW (cheW), CheX (cheX) and CheY (cheY) genes, complete cds	98	397
63	13	5012	4779 gblU128611	Treponema pallidum Nichols TpN38(b) (tpn38(b)) gene, complete cds	100	186
63	14	5484	4945 gblU128611	Treponema pallidum Nichols TpN38(b) (tpn38(b)) gene, complete cds	100	458

TABLE 1.

Treponema pallidum - Coding regions containing known sequences

63	15	5648	5409	gbIU12861I	Treponema pallidum Nichols TpN38(b) (TpN38(b))	100	190
63	22	11392	10127	gbIU02628I	Treponema pallidum pallidum Nichols TpN30 precursor (TpN50) gene, complete cds	99	1266
63	23	12291	11371	gbIU28427I	Treponema pallidum antigen (tpp57) gene, complete cds	99	312
64	9	3849	5198	gbIU04241I	T. pallidum 34 kd antigen gene, complete cds	99	177
64	10	5102	5317	gbIU04241I	T. pallidum 34 kd antigen gene, complete cds	100	216
64	11	5364	5795	gbIU04241I	T. pallidum 34 kd antigen gene, complete cds	96	382
64	12	6257	5658	gbIU04241I	T. pallidum 34 kd antigen gene, complete cds	93	414
64	13	5726	5965	gbIU04241I	T. pallidum 34 kd antigen gene, complete cds	100	200
64	14	6082	7029	gbIU04241I	T. pallidum 34 kd antigen gene, complete cds	96	404
68	1	512	3	gbIU88957I	Treponema pallidum major outer sheath protein homolog Msp (msp) gene, complete cds	97	482
69	6	2742	3845	gbIU61535I	Treponema pallidum gamma-glutamyl phosphate reductase (proA) and glutamate 5-kinase (prob) genes, complete cds	98	791
69	7	3578	4021	gbIU61535I	Treponema pallidum gamma-glutamyl phosphate reductase (proA) and glutamate 5-kinase (prob) genes, complete cds	99	428
69	8	3994	4596	gbIU61535I	Treponema pallidum gamma-glutamyl phosphate reductase (proA) and glutamate 5-kinase (prob) genes, complete cds	99	546
69	9	4491	4907	gbIU61535I	Treponema pallidum gamma-glutamyl phosphate reductase (proA) and glutamate 5-kinase (prob) genes, complete cds	100	370
69	10	4949	5296	gbIU61535I	Treponema pallidum gamma-glutamyl phosphate reductase (proA) and glutamate 5-kinase (prob) genes, complete cds	100	52
72	2	1313	312	gbIU97360I	Treponema pallidum HFLK homolog gene, complete cds	99	609
72	3	1137	2792	gbIU97360I	Treponema pallidum HFLK homolog gene,	99	553

TABLE 1.

Treponema pallidum - Coding regions containing known sequences

78	1	790	2	gb U36839	Treponema pallidum putative switch protein FliM (fliM) gene, partial cds, and FliY (fliY) gene, complete cds, export apparatus proteins FliP (fliP), FliQ (fliQ), FliR (fliR) and FliB (fliB), signal transducing receptor FliA (fliA), GTP binding protein >	complete cds	99	716
78	2	1916	1074	gb U36839	Treponema pallidum putative switch protein FliM (fliM) gene, partial cds, and FliY (fliY) gene, complete cds, export apparatus proteins FliP (fliP), FliQ (fliQ), FliR (fliR) and FliB (fliB), signal transducing receptor FliA (fliA), GTP binding protein >	complete cds	100	843
78	3	2320	1886	gb U36839	Treponema pallidum putative switch protein FliM (fliM) gene, partial cds, and FliY (fliY) gene, complete cds, export apparatus proteins FliP (fliP), FliQ (fliQ), FliR (fliR) and FliB (fliB), signal transducing receptor FliA (fliA), GTP binding protein >	complete cds	100	373
78	4	2703	2224	gb U36839	Treponema pallidum putative switch protein FliM (fliM) gene, partial cds, and FliY (fliY) gene, complete cds, export apparatus proteins FliP (fliP), FliQ (fliQ), FliR (fliR) and FliB (fliB), signal transducing receptor FliA (fliA), GTP binding protein >	complete cds	99	445
78	5	3709	2609	gb U36839	Treponema pallidum putative switch protein FliM (fliM) gene, partial cds, and FliY (fliY) gene, complete cds, export apparatus proteins FliP (fliP), FliQ (fliQ), FliR (fliR) and FliB (fliB), signal transducing receptor FliA (fliA), GTP binding protein >	complete cds	99	721
78	6	3998	3645	gb U28219	Treponema pallidum flagellar hook (fliE), (orf4), flagellar motor protein (motA), flagellar motor	complete cds	98	300

TABLE 1.

Treponema pallidum - Coding regions containing known sequences

78	7	4483	3707	gbIU28219I	Treponema pallidum flagellar hook (fliE), (orf4), flagellar motor protein (motA), flagellar motor protein (motB), (fliL), and flagellar switch protein (fliM) genes, complete cds, and (fliY) gene, partial cds	100	388	
78	8	4500	4243	gbIU28219I	Treponema pallidum flagellar hook (fliE), (orf4), flagellar motor protein (motA), flagellar motor protein (motB), (fliL), and flagellar switch protein (fliM) genes, complete cds, and (fliY) gene, partial cds	100	126	
81	6	3472	2957	gbIM17716I	T.pallidum basic membrane protein gene, complete cds	99	495	
81	7	4116	3268	gbIM17716I	T.pallidum basic membrane protein gene, complete cds	100	612	
81	8	6081	4114	gbIM17716I	T.pallidum basic membrane protein gene, complete cds	100	211	
83	4	1058	1795	embIX61227I	T.pallidum Tp70 gene (partial)	98	179	
83	14	5091	5687	gbIU97362I	Treponema pallidum 22.5K protein gene, complete cds	100	80	
83	15	5828	6454	gbIU97362I	Treponema pallidum 22.5K protein gene, complete cds	100	627	
83	16	6617	7972	gbIU97362I	Treponema pallidum 22.5K protein gene, complete cds	97	77	
84	1	3	152	gbIU61851I	Treponema pallidum histidine kinase CheA (cheA), and chemotaxis proteins CheW (cheW), CheX (cheX) and CheY (cheY) genes, complete cds	95	96	
84	2	659	393	gbIU61851I	Treponema pallidum histidine kinase CheA (cheA), and chemotaxis proteins CheW (cheW), CheX (cheX) and CheY (cheY) genes, complete cds	100	106	

TABLE 1.

Treponema pallidum - Coding regions containing known sequences

84	3	837	598	gbIU708681	Treponema pallidum ribosomal protein L28 homolog (rpL28) gene, complete cds	86	181
84	4	703	1236	gbIU708681	Treponema pallidum ribosomal protein L28 homolog (rpL28) gene, complete cds	88	224
89	6	2962	4770	gbIU577571	Treponema pallidum PolA gene, partial cds	99	330
89	7	4849	5748	gbIU577571	Treponema pallidum PolA gene, partial cds	100	837
89	8	5678	8002	gbIU577571	Treponema pallidum PolA gene, partial cds	99	2325
89	9	8378	8034	gbIU577571	Treponema pallidum PolA gene, partial cds	99	201
99	1	87	1316	gbIL203011	Treponema pallidum 38 kilodalton glucose/galactose binding lipoprotein (lpp38) gene, complete cds	99	1230
99	2	1379	2926	gbIU484161	Treponema pallidum mgl operon: putative glucose/galactose binding protein (mglB), putative ATP binding protein (mglA), and hydrophobic putative membrane-associated protein (mglC) genes, complete cds	99	1548
99	3	2806	3483	gbIU453231	Treponema pallidum putative ATP-binding protein (mglA) and putative methylgalactoside transport protein (mglC) genes, complete cds	99	624
99	4	3306	3043	gbIU453231	Treponema pallidum putative ATP-binding protein (mglA) and putative methylgalactoside transport protein (mglC) genes, complete cds	98	264
99	5	3305	4534	gbIU484161	Treponema pallidum mgl operon: putative glucose/galactose binding protein (mglB), putative ATP binding protein (mglA), and hydrophobic putative membrane-associated protein (mglC) genes, complete cds	99	1113
99	6	6086	4920	gbIU706611	Treponema pallidum GTP-binding protein, mccF-like protein and ATP-dependent DNA helicase (RecG) genes, complete cds	99	1101
100	31	12458	13066	emblX639651	T. pallidum endoflagellar genes ftaB1 and ftaB3	99	566
100	32	13027	13410	emblX639651	T. pallidum endoflagellar genes ftaB1 and ftaB3	100	370

TABLE 1.

Treponema pallidum - Coding regions containing known sequences

			TPFLAB13		
100	33	13658	13957 embIX63965I	T.pallidum endoflagellar genes f1ab1 and flaB3	100 262
100	34	13846	14544 embIX63965I	T.pallidum endoflagellar genes f1ab1 and flaB3	98 628
104	7	2890	2552 gb M74825I	Treponema pallidum 17 kDa lipoprotein gene, complete cds	99 256
109	1	150	8666 gb U88957I	Treponema pallidum major outer sheath protein homolog Msp (msp) gene, complete cds	100 143
112	5	3982	3023 gb U56999I	Treponema pallidum methyl-accepting chemotaxis protein (mcp-1) gene, complete cds, and potential regulatory molecule (pfos/R) gene, partial cds	100 791
112	6	5935	4076 gb U56999I	Treponema pallidum methyl-accepting chemotaxis protein (mcp-1) gene, complete cds, and potential regulatory molecule (pfos/R) gene, partial cds	100 1818
112	7	6555	5866 gb U56999I	Treponema pallidum methyl-accepting chemotaxis protein (mcp-1) gene, complete cds, and potential regulatory molecule (pfos/R) gene, partial cds	99 477
112	8	7024	6545 gb U56999I	Treponema pallidum methyl-accepting chemotaxis protein (mcp-1) gene, complete cds, and potential regulatory molecule (pfos/R) gene, partial cds	100 480
112	9	6572	7804 gb U56999I	Treponema pallidum methyl-accepting chemotaxis protein (mcp-1) gene, complete cds, and potential regulatory molecule (pfos/R) gene, partial cds	100 939
112	10	7282	7022 gb U56999I	Treponema pallidum methyl-accepting chemotaxis protein (mcp-1) gene, complete cds, and potential regulatory molecule (pfos/R) gene, partial cds	100 261
112	11	7819	7280 gb U56999I	Treponema pallidum methyl-accepting chemotaxis protein (mcp-1) gene, complete cds, and potential regulatory molecule (pfos/R) gene, partial cds	100 231
112	17	13968	13516 embIX5411II	T.pallidum Tp4 gene	99 449
120	1	850	44 gb U55214I	Treponema pallidum RuvB' gene, partial cds, and	100 779

TABLE 1.

Treponema pallidum - Coding regions containing known sequences

			TroA, TroB, TroC, TroD, TroR, Pgm, Pf, and Tpp15 genes, complete cds
121	23	12559	13242 gblM1093111 T.pallidum tnpA gene encoding a membrane protein
121	24	13291	14418 gblM1093111 T.pallidum tnpA gene encoding a membrane protein
121	25	14361	14804 gblM5856211 Treponema pallidum tnpB gene, complete cds
121	26	14750	15178 gblM5856211 Treponema pallidum tnpB gene, complete cds
132	1	2	724 gblM8876911 Treponema pallidum 47-kilodalton antigen gene, complete cds
133	1	709	2 gblU7374811 Treponema pallidum putative aspartate aminotransferase TpAAAT (tpaat) and leucine-rich repeat protein TpLRR genes, complete cds
133	2	1346	888 gblU7374811 Treponema pallidum putative aspartate aminotransferase TpAAAT (tpaat) and leucine-rich repeat protein TpLRR genes, complete cds
133	3	2265	1549 gblU7374811 Treponema pallidum putative aspartate aminotransferase TpAAAT (tpaat) and leucine-rich repeat protein TpLRR genes, complete cds
142	3	4022	4951 gblU7809411 Treponema pallidum DNA gyrase subunit B (gyrB) gene, complete cds
142	4	5107	6972 gblU7809411 Treponema pallidum DNA gyrase subunit B (gyrB) gene, complete cds
145	12	3527	3958 gblU8895711 Homolog Msp (msp) gene, complete cds
153	2	1801	746 gblU3268311 Treponema pallidum major outer sheath protein A (cfpA) gene, complete cds
166	1	19	843 gblU5521411 Treponema pallidum RuvB' gene, partial cds, and TroA, TroB, TroC, TroD, TroR, Pgm, Pf, and Tpp15 genes, complete cds
166	2	689	994 gblU5521411 Treponema pallidum RuvB' gene, partial cds, and TroA, TroB, TroC, TroD, TroR, Pgm, Pf, and

TABLE 1.

Treponema pallidum - Coding regions containing known sequences

			Tpp15 genes, complete cds		
166	3	943	1710 gblU552141	Treponema pallidum RuvB' gene, partial cds, and TroA, TroB, TroC, TroD, TroR, Pgm, PfS, and Tpp15 genes, complete cds	99 750
166	4	1708	2655 gblU552141	Treponema pallidum RuvB' gene, partial cds, and TroA, TroB, TroC, TroD, TroR, Pgm, PfS, and Tpp15 genes, complete cds	98 915
166	5	2594	2821 gblU552141	Treponema pallidum RuvB' gene, partial cds, and TroA, TroB, TroC, TroD, TroR, Pgm, PfS, and Tpp15 genes, complete cds	98 164
166	6	2814	3287 gblU552141	Treponema pallidum RuvB' gene, partial cds, and TroA, TroB, TroC, TroD, TroR, Pgm, PfS, and Tpp15 genes, complete cds	100 474
166	7	3352	3831 gblU552141	Treponema pallidum RuvB' gene, partial cds, and TroA, TroB, TroC, TroD, TroR, Pgm, PfS, and Tpp15 genes, complete cds	100 456
166	8	3714	4061 gblU552141	Treponema pallidum RuvB' gene, partial cds, and TroA, TroB, TroC, TroD, TroR, Pgm, PfS, and Tpp15 genes, complete cds	99 254
168	1	1113	4 gblU368391	Treponema pallidum putative switch protein FilM (filM) gene, partial cds, and FilY (filY) gene, complete cds, export apparatus proteins FilP (filP), FilQ (filQ), FilR (filR) and FilB (filB), signal transducing receptor FlhA (flhA), GTP binding protein >	99 1043
172	3	1442	498 gblU889571	Treponema pallidum major outer sheath protein homolog Msp (msp) gene, complete cds	97 98
172	4	519	989 gblU889571	Treponema pallidum major outer sheath protein homolog Msp (msp) gene, complete cds	84 119
174	2	1640	372 gblU889571	Treponema pallidum major outer sheath protein homolog Msp (msp) gene, complete cds	90 845
176	1	71	835 gblU282191	Treponema pallidum flagellar hook (flgE), (orf4), flagellar motor protein (motA), flagellar motor	99 681

TABLE 1.

Treponema pallidum - Coding regions containing known sequences

			protein (motB), (fliL), and flagellar switch protein (fliM) genes, complete cds, and (fliY) gene, partial cds		
176	2	899	1228 gbuU282191 Treponema pallidum flagellar hook (fliE), (orf4), flagellar motor protein (motA), flagellar motor protein (motB), (fliL), and flagellar switch protein (fliM) genes, complete cds, and (fliY) gene, partial cds	100	219
190	1	3	383 gbuU618511 Treponema pallidum histidine kinase CheA (cheA), and chemotaxis proteins CheW (cheW), CheX (cheX) and CheY (cheY) genes, complete cds	98	211
190	2	248	754 gbuU618511 Treponema pallidum histidine kinase CheA (cheA), and chemotaxis proteins CheW (cheW), CheX (cheX) and CheY (cheY) genes, complete cds	99	441
205	1	2	568 gbuU657431 Treponema pallidum outer membrane protein (trmp2) gene, complete cds	98	565
206	2	528	1214 gbuU552141 Treponema pallidum RuvB' gene, partial cds, and TroA, TroB, TroC, TroD, TroR, Pgm, Pf, and Tpp15 genes, complete cds	98	177
206	3	1183	740 gbuU552141 Treponema pallidum RuvB' gene, partial cds, and TroA, TroB, TroC, TroD, TroR, Pgm, Pf, and Tpp15 genes, complete cds	98	146
206	4	1025	804 gbuU552141 Treponema pallidum RuvB' gene, partial cds, and TroA, TroB, TroC, TroD, TroR, Pgm, Pf, and Tpp15 genes, complete cds	98	63
206	5	1599	1204 gbuU552141 Treponema pallidum RuvB' gene, partial cds, and TroA, TroB, TroC, TroD, TroR, Pgm, Pf, and Tpp15 genes, complete cds	99	396
223	1	1	564 embBX541111 T pallidum Tp4 gene TP TP4	99	358
224	3	992	453 gbuU889571 Treponema pallidum major outer sheath protein homolog Msp (msp) gene, complete cds	97	411
224	4	486	731 gbuU889571 Treponema pallidum major outer sheath protein	100	154

TABLE 1.

Treponema pallidum - Coding regions containing known sequences

				Treponema pallidum - Coding regions containing known sequences
231	1	447	46	gbIU889571 [homolog Msp (msp) gene, complete cds
491	1	394	83	gbIU952141 [Treponema pallidum major outer sheath protein homolog Msp (msp) gene, complete cds
495	1	406	2	gbIU618511 [Treponema pallidum histidine kinase CheA (cheA), synthase gene, complete cds
495	2	582	235	gbIU618511 [Treponema pallidum histidine kinase CheA (cheA), and chemotaxis proteins CheW (cheW), CheX (cheX) and CheY (cheY) genes, complete cds
562	1	718	44	gbIU889571 [Treponema pallidum major outer sheath protein (cheX) and CheY (cheY) genes, complete cds
587	1	110	439	gbIU889571 [Treponema pallidum major outer sheath protein homolog Msp (msp) gene, complete cds
592	1	3	410	gbIU042411 [T. pallidum 34 kd antigen gene, complete cds
602	1	1	246	gbIU569991 [Treponema pallidum methyl-accepting chemotaxis protein (mcp-1) gene, complete cds, and potential regulatory molecule (proS/R) gene, partial cds
626	2	642	73	gbIU889571 [Treponema pallidum major outer sheath protein homolog Msp (msp) gene, complete cds
627	1	402	88	gbIU889571 [Treponema pallidum major outer sheath protein homolog Msp (msp) gene, complete cds
627	2	911	366	gbIU889571 [Treponema pallidum major outer sheath protein homolog Msp (msp) gene, complete cds
630	1	1	549	gbIU569991 [Treponema pallidum methyl-accepting chemotaxis protein (mcp-1) gene, complete cds, and potential regulatory molecule (proS/R) gene, partial cds
630	2	165	4	gbIU569991 [Treponema pallidum methyl-accepting chemotaxis protein (mcp-1) gene, complete cds, and potential regulatory molecule (proS/R) gene, partial cds
633	2	102	332	gbIU889571 [Treponema pallidum major outer sheath protein

TABLE 1. Treponema pallidum - Coding regions containing known sequences

				homolog Msp (msp) gene, complete cds		
640	1	708	310	gbIU889571 Treponema pallidum major outer sheath protein	93	.274
645	1	3	311	gbIU889571 homolog Msp (msp) gene, complete cds	85	122
646	3	896	501	gbIU889571 Treponema pallidum major outer sheath protein	91	148
653	1	1	282	emblX639651 T.pallidum endoflagellar genes flAB1 and flAB3	97	126
653	2	499	161	embIX639651 T.pallidum endoflagellar genes flAB1 and flAB3	98	272
654	1	1	552	gbIU889571 Treponema pallidum major outer sheath protein	68	293
654	2	624	4	gbIU889571 Treponema pallidum major outer sheath protein	68	293
660	1	585	4	gbIU889571 homolog Msp (msp) gene, complete cds	71	252
665	1	434	3	gbIU889571 Treponema pallidum major outer sheath protein	84	119
665	2	57	434	gbIU889571 homolog Msp (msp) gene, complete cds	84	119
666	1	2	373	gbIU973581 Treponema pallidum 29K protein gene, complete	98	178
666	2	496	119	gbIU973581 Treponema pallidum 29K protein gene, complete	98	301
669	2	412	71	gbIU889571 Treponema pallidum major outer sheath protein	75	203
671	1	361	62	gbIU889571 homolog Msp (msp) gene, complete cds	100	297
672	1	3	653	gbIU889571 Treponema pallidum major outer sheath protein	90	419
681	1	109	258	gbIU780941 Treponema pallidum DNA gyrase subunit B (gyrB)	94	118

TABLE 1.

Treponema pallidum - Coding regions containing known sequences

696	1	304	483	gbIU453231	Treponema pallidum putative ATP-binding protein (mgA) and putative methylgalactoside transport protein (mgC) genes, complete cds	97	129
-----	---	-----	-----	------------	--	----	-----

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident
40	1	340	316	gpiX57972	beta-lactamase TEM6 [Escherichia coli] >pirSIS24415 beta-lactamase [EC 3.5.2.6] TEM6 - Escherichia coli	100	100
40	2	965	405	gpiM74750	beta-lactamase [Escherichia coli] >gpiD14640 ECORPL12A_2 beta-lactamase [Plasmid pKF2] >gpiU36911 SAU36911_1 beta-lactamase [Staphylococcus aureus] >gpiU36912 SAU36912_1 beta-lactamase [Staphylococcus aureus]	100	99
61	1	418	972	gpiM74750	beta-lactamase [Escherichia coli] >gpiD14640 ECORPL12A_2 beta-lactamase [Plasmid pKF2] >gpiU36911 SAU36911_1 beta-lactamase [Staphylococcus aureus] >gpiU36912 SAU36912_1 beta-lactamase [Staphylococcus aureus]	100	98
96	1	12	275	pirSIA4356_2	>gpiU25060 XXU25060_3 beta lactamase [unident] homeotic protein Hox 4.3 - mouse	100	85
228	2	832	1425	gpiM74750	beta-lactamase [Escherichia coli] >gpiD14640 ECORPL12A_2 beta-lactamase [Plasmid pKF2] >gpiU36911 SAU36911_1 beta-lactamase [Staphylococcus aureus] >gpiU36912 SAU36912_1 beta-lactamase [Staphylococcus aureus]	100	100
534	1	107	616	gpiJ02459	>gpiSIVHBPEL major capsid protein E - phage lambda Bacteriophage lambda, complete genome. [Bacteriophage lambda]	100	100
579	1	2	241	gpiX66453	>pirSIVNCSPBL_1 DNA-directed RNA polymerase [Euploites octocarinatus]	100	87
637	1	29	394	gpiA04190	galactosidase fusion protein [unidentified]	100	99
					>gpIK01075 SYNCSPBL_1 Plasmodium knowlesi circumsporozoite protein repeat region fused with beta-lactamase gene of pBR322. [Artificial gene] [SUB 181-300]		
					>gpIL30112P14BLAREP_1 beta-lactamase [Plasmid pR01614]		
45	16	15198	14485	pirSIA3705_3	33K endoflagellar protein FlabB - Treponema pallidum 33K class B flagellar protein, periplasmic - Treponema pallidum subsp. pallidum (fragment) [SUB 1-21]	98	98
607	1	223	2	gpiU13869	lacZ alpha peptide [Cloning vector]	98	98

TABLE 2.
Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

642	1	81	392	gpiM74750	beta-lactamase [Escherichia coli] >gpiD14640[ECORPL12A_2 beta-lactamase [Plasmid pKF2] >gpiU36911[ISAU3691_1 beta-lactamase [Staphylococcus aureus] >gpiU36912[ISAU36912_1 beta-lactamase [Staphylococcus aureus]	98	98
14	1	3	329	gpiU250591	LacZ alpha peptide [unidentified] >gpiU25060[XXU25060_2 LacZ alpha peptide [unidentified] >gpiU25061[XXU25061_2 LacZ alpha peptide [unidentified] >gpiU23751[CVU23751_2 beta galactosidase [Cloning vector pBBR1MCS-2]	96	96
101	12	5676	5107	gpiZ677531	DNA-replication helicase [Chloroplast Odontella sinensis]	96	70
169	1	2	190	gpiV000831	Artificial cloning vehicle pBR327, derived from pBR322. The sequence was not resequenced but deduced from the pBR322 sequence. Contains the reading frames for ampicillin resistance (Apr) and tetracycline resistance (Tcr) and an origin of replication. [uni]	96	96
588	1	1	213	gpiX946071	MefI protein [Saccharomyces cerevisiae]	96	74
11	15	6934	8265	gpiX045811	E.coli recB gene for exonuclease V. [Escherichia coli] >gpiU2958[IECU2958_31 exonuclease V subunit [Escherichia coli] >gpiU2958[IECU2958_31 exonuclease V subunit [Escherichia coli] >pirSINCECX5 exodeoxyribonuclease V (EC 3.1.11.5) 135K chain - Escherichia coli	95	65
34	1	328	2	gpiU374561	beta-lactamase [Cloning vector YITAG100] >gpiU37457[CVU37457_2 beta-lactamase [Cloning vector YATAG200]	95	94
673	1	376	2	gpiA041901	galactosidase fusion protein [unidentified]	95	95
					>gpiK1075[SYNCSPBL_1 Plasmodium knowlesi circumsporozoite protein repeat region fused with beta-lactamase gene of pBR322. [Artificial gene] [SUB 181-300]		
					>gpiL30112P14BLAREP_1 beta-lactamase [Plasmid pR01614]		
79	9	4159	5235	gpiU156091	flagellar switch protein [Treponema denticola]	94	85
117	2	1565	921	gpiD640061	>gpiL3685[ITRPFLIG_1 flIG gene product [Treponema denticola]	93	76
16	25	15399	16082	gpiL103281	hypothetical protein [Synechocystis sp.]	92	71

TABLE 2. *Treponema pallidum* - Putative coding regions of novel proteins similar to known proteins

				>gpiX0163 EC0UNC_3 E. coli origin of replication oriC and genes gnd, unc, EcoURF-1 and glmS. [Escherichia coli]	
				>priSIBVECQB gndB protein - Escherichia coli	
42	27	18334	18720	gpiMS7776 ribosomal protein S12 [Leptospira biflexa] >priSIA36152	91
				ribosomal protein S12 - Leptospira biflexa (serotype patoc)	82
197	1	393	370	gpiID640001 hypothetical protein [Synechocystis sp.]	91
547	1	552	370	gpiUT9381 unknown [Klebsiella pneumoniae]	81
16	49	29062	28667	gpiM244661 S.typhimurium flagellar L-ring (flgH), flagellar P-ring (flgI), and flagellar (flgJ) genes, complete cds. [Salmonella typhimurium] >priSIC30930 flagellar basal body protein flgJ - Salmonella typhimurium	90
650	1	179	391	gpiUT171391 Phagemid cloning vector pSIT, complete sequence. [Cloning vector pSIT] >gpiU47102 CVU47102_2 beta-lactamase [Cloning vector pALTER<R>-Ex1]	89
4	8	6487	5276	gpiX52898 glyceraldehyde 3-phosphate dehydrogenase [Trypanosoma cruzi] >gpiX52898 TCGAP_2 glyceraldehyde 3-phosphate dehydrogenase [Trypanosoma cruzi] >priSDEUT1C glyceraldehyde-3-phosphate dehydrogenase (EC 1.2.1.12), glycosomal - Trypanosoma cruzi	89
8	45	26587	27414	gpiD306901 HSP40 [Staphylococcus aureus]	72
25	1	224	3	priSIA6016 glycoprotein IIb - rat	89
64	8	3423	3635	gpiZ190591 Cek8 protein [Gallus gallus] >gpiZ19059 GGCEK8A_1 Cek8 protein [Gallus gallus] >priSIS33505 protein-tyrosine kinase (EC 2.7.1.112) Cek8 - chicken (fragment) [SUB 2-849] >gpiX5724 IMMMMPK3_1 tyrosine kinase [Mus musculus] [SUB 611-668] >priSIP10183 protein	68
598	1	361	212	gpiL45070 glutamine amidotransferase [Haemophilus influenzae] >gpiU32726 HU32726_10 glutamine amidotransferase [Haemophilus influenzae] >gpiU00073 HU00073_51 glutamine amidotransferase [Haemophilus influenzae] >gpiU32835 HU32835_5 glutamine amidotransferase [Hae	72

TABLE 2.
Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

8	17	8697	9533 gpiI36380	tufA gene product [Neisseria gonorrhoeae]	87	76
16	27	16417	16127 gpiFI4628	cyclophilin B [Sus scrofa]	87	79
30	9	4454	4765 gpiU13165	sigma 43 subunit of RNA polymerase [Listeria monocytogenes]	87	50
44	16	9669	9334 gpiI45205	asparagine synthetase A [Haemophilus influenzae] >gpiU32738 HTU32738_3 asparagine synthetase A [Haemophilus influenzae] >gpiU0074 HTU0074_84 asparagine synthetase A [Haemophilus influenzae] >gpiU32847 HTU32847_3 asparagine synthetase [Haemophilus influenzae]	87	76
143	2	950	768 gpiS63735	membrane glycoprotein M6 [Mus sp.]	87	62
42	22	15120	15563 gpiM38305	E.coli RNA polymerase beta subunit (rpoC) gene, partial cds. [Escherichia coli]	86	72
48	18	11952	11773 gpiI45145	periplasmic ribose-binding protein [Haemophilus influenzae] >gpiU32732 HTU32732_5 periplasmic ribose-binding protein [Haemophilus influenzae] >gpiU00074 HTU00074_25 periplasmic ribose-binding protein [Haemophilus influenzae] >gpiU32841 HTU32841_5 D-ribose	86	50
113	2	1687	2058 gpiZ9124	SecA [Chloroplast Spinacia olereacea]	86	65
113	4	2604	2900 gpiX55034	SecA protein [Escherichia coli] >gpiD 0483 ECO110K_78 secA protein [Escherichia coli] >gpm2079 IECOSECA_2 secA gene product [Escherichia coli] >pirSIA31088 secA protein - Escherichia coli >gpm1921 IECOENVAA_4 secA gene product [Escherichia coli] [SUB 1]	86	65
117	3	2176	1472 gpiI45353	ATP-dependent clip protease proteolytic component [Haemophilus influenzae] >gpiU32754 HTU32754_5 ATP-dependent clip protease proteolytic component [Haemophilus influenzae] >gpiU0076 HTU0076_53 ATP-dependent clip protease proteolytic component [Haemophilus elongation factor Tu [Escherichia coli]	86	65
8	19	9481	10374 gpiJ01690	>gpiU18997 ECOUW67_263 tufA gene product [Escherichia coli] >pirSIEFFECTA translation elongation factor Tu.A - Escherichia coli	85	68
63	12	4803	4378 gpiU14003	Escherichia coli K-12 chromosomal region from 92.8 to 00.1	85	74

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

588	2	557	117	gpLI4925I	minutes. [Escherichia coli] >priSIS5631S5 hypothetical protein 510		
693	1	64	441	gpIM96436I	- Escherichia coli >gpU00061ECOUW89_134 E. coli		
					chromosomal region from 89.2 to 92.8 minutes. [Escherichia coli]		
					{SUB 9-510}		
					nitrogen fixation protein [Anabaena sp.]		
					oxaloacetate decarboxylase [Klebsiella pneumoniae] >priSIA44464		
					oxaloacetate decarboxylase beta subunit (C terminus) - Klebsiella		
					pneumoniae (fragment) >gpIM26290IKPN0ADBG2_1 oadB gene		
					product [Klebsiella pneumoniae] {SUB 4-141}		
1	5	1993	2015	gpIL44957I	Holliday junction DNA helicase [Haemophilus influenzae]		
					>gpU32716HTU32716_13 Holliday junction DNA helicase		
					[Haemophilus influenzae] >gpU00072HTU00072_42 Holliday		
					junction DNA helicase [Haemophilus influenzae]		
					>gpU32825HTU32825_12 Holliday junction		
21	20	11250	11630	gpIZ46729I	unknown [Saccharomyces cerevisiae] >priSIS49802 hypothetical		
					protein YM9958.03c - yeast (Saccharomyces cerevisiae)		
47	25	14251	13676	priSIC4715	ribosomal protein S16 - Bacillus subtilis		
			4				
49	1	1	1092	gpIX8941II	SecA [Rhodobacter capsulatus]		
79	11	6479	7015	gpIM72718I	B. subtilis ftaA locus operon. [Bacillus subtilis]		
					>gpX56049IBSFLAAO_3_B - subtilis ftaA locus operon.		
					[Bacillus subtilis] >priSIPWB5AS H+-transporting ATP synthase		
					alpha chain homolog - Bacillus subtilis		
87	2	1995	1639	gpIZ5417II	elongation factor EF-G [Rickettsia prowazekii]		
					>gpU02603RP0U02603_8 elongation factor EF-G [Rickettsia		
					prowazekii] {SUB 1-47}		
101	23	15080	13362	gpIU43536I	clpB gene product [Corynebacterium glutamicum]		
8	43	24498	26312	gpIX67646I	heat-shock protein [Borrelia burgdorferi]		
					>gpIM96847IBORGPEPLS_2_dnaK homologue gene product		
					[Borrelia burgdorferi] >gpIM97912BORHSP70A_170 kDa heat		
					shock protein [Borrelia burgdorferi] >gpIS42385IS42385_1		
					HSP70 homolog [Borrelia burgdorferi] CA12 isol		

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

10	21	12836	11373] gpiU35673]	Rho [Borrelia burgdorferi] >pirSIS35618 rho protein - Lyme disease spirochete (fragment) {SUB 97-515} >gpl076561BORRHO_1 Rho protein [Borrelia burgdorferi] {SUB 127-515}	83	64
23	6	2510	3373] gpiM64730]	DNA mismatch repair protein [Escherichia coli] >gplU29579IECU29579_30 DNA mismatch repair protein [Escherichia coli]	83	66
45	1	951	4 gpiM77351]	ATP-binding protein [Streptococcus mutans] >pirSIE42400 ATP-binding protein MsmK - Streptococcus mutans >pirSIC27626 hypothetical protein 2 - Streptococcus mutans (fragment) {SUB 1-33}	83	67
61	14	8382	7999 gpiU454261]	heat shock protein HTPG [Borrelia burgdorferi] >gplL32145IBORHHTPG_1 C62.5 heat shock protein [Borrelia burgdorferi] {SUB 497-5751}	83	56
81	5	2415	2224] gpiX908571]	-14 gene product [Homo sapiens] IepA gene product [Bacillus subtilis] >gpiD17650IBACGPR_4	83	50
115	6	2113	2598 gpiX916351]	ORF80 protein [Bacillus subtilis] {SUB 1-327} ORF3 [Bacillus subtilis] >pirSUN0146 hypothetical protein (div+ 3' region) - Bacillus subtilis (fragment)	83	61
41	9	2379	3092 gpiD102791]	groES gene product [Bacillus stearothermophilus] >pirSIA49855 heat shock protein GroES - Bacillus stearothermophilus permease [Haemophilus influenzae]	82	67
81	1	287	628 gpiL101321]	groES gene product [Bacillus stearothermophilus] >pirSIA49855 heat shock protein GroES - Bacillus stearothermophilus permease [Haemophilus influenzae]	82	60
165	6	1887	2087 gpiU32690]	ribosomal protein S10 [Haemophilus influenzae] >gpiU32761IHU32761_6 ribosomal protein S10 [Haemophilus influenzae] >gpiU00077IHU00077_35 ribosomal protein S10 [Haemophilus influenzae] >gpiU32707IHU32707_9 30S ribosomal protein S10 [Haemophilus influenzae] VP1 [Unknown.]	82	54
8	20	10445	10771 gpiL454141]	asparagine synthetase A [Haemophilus influenzae] >gpiU32738IHU32738_3 asparagine synthetase A [Haemophilus influenzae] >gpiU00074IHU00074_84 asparagine synthetase A [Haemophilus influenzae] >gpiU32847IHU32847_3 asparagine	81	70
44	17	10376	9630 gpiL45205]		81	54

TABLE 2. Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

47	1	900	4	gpi D49951	H2O-forming NADH Oxidase [Streptococcus mutans]	81	64
54	28	8744	8499	gpi U396634	purine-nucleoside phosphorylase [Mycoplasma genitalium] >pir SISD64205 purine-nucleoside phosphorylase - Mycoplasma genitalium (SGC3)	81	53
57	11	7418	9556	gpi L77216	Lon protease [Borrelia burgdorferi]	81	59
57	12	9477	9857	gpi 03896	E.coli ATP-dependent protease La (lon) gene, complete cds. [Escherichia coli]	81	66
101	25	16036	14795	gpi U43536	cipB gene product [Corynebacterium glutamicum]	81	69
112	1	912	4	gpi Z12160	gidA gene product [Borrelia burgdorferi] >gpi Z12160 BBGIDAG_1 division protein [Borrelia burgdorferi] (SUB 529-593) >gpi X95669BBTHDFGID_2 gidA gene product [Borrelia burgdorferi] (SUB 1-29) >gpi X95668BBGIDMOXR_1	81	64
174	1	117	317	gpi X15540	gidA gene product [Borrelia burgdorferi] T.brucei pg1 gene for glucose-6-phosphate isomerase (EC 5.3.1.9). [Trypanosoma (Trypanozoon) brucei] >pir SINUUUTB glucose-6-phosphate isomerase (EC 5.3.1.9) - Trypanosoma brucei	81	63
337	1	84	374	gpi M25927	Influenza A/swine/Hong Kong/1/26/82 (H3N2) PB1 gene, complete cds. [Influenza virus type A]	81	54
555	1	2	190	gpi L14925	nitrogen fixation protein [Anabaena sp.]	81	79
683	1	388	170	gpi X69618	inhibin alpha-subunit [Mus musculus] >pir SIS31439 inhibin alpha chain - mouse	81	72
8	28	16269	16844	pir SIR5BS5	ribosomal protein L5 - Bacillus stearothermophilus	80	57
16	16	11479	10670	gpi X77515	F pyruvate oxidoreductase [Rhodospirillum rubrum] >pir SIS41961	80	67
17	9	4729	4893	gpi L11706	pyruvate oxidoreductase - Rhodospirillum rubrum hormone-sensitive lipase [Homo sapiens]	80	80
35	16	9193	9642	gpi L04500	thioredoxin reductase [Eubacterium acidaminophilum] >pir SIS36988 thioredoxin reductase chain B - Eubacterium acidaminophilum	80	59
36	16	10409	10756	gpi X78993	probable transfer RNA-Gly synthetase [Saccharomyces cerevisiae]	80	55

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				>gplZ3590[SCYBR12]C_1 GSR1 gene product [Saccharomyces cerevisiae] >pirSIS48285 probable glycine- <u>tRNA</u> ligase (EC 6.1.1.14) GRS1 - yeast (Saccharomyces cerevisiae)		
42	19	11240	13345	gplI48488I RNA polymerase beta subunit [Borrelia burgdorferi]	80	67
42	20	13225	13983	gplI48488I RNA polymerase beta subunit [Borrelia burgdorferi]	80	72
42	21	13981	15210	gplI43593I RNA polymerase beta' subunit [Bacillus subtilis] >gplI43593IBACBPSO_1 RNA polymerase beta' subunit [Bacillus subtilis] >gplI24376BACRPLL_4 RNA polymerase beta-subunit [Bacillus subtilis] (SUB 1-466)	80	65
45	3	2471	5083	gplM31045I E.coli ATP-dependent Clp protease (clpA) gene, complete cds. [Escherichia coli] >pirSISUECCA ATP-dependent Clp protease (EC 3.4.21.-) chain A - Escherichia coli >gplM23220IECCOCLPA_1 ATP-dependent protease [Escherichia coli] (SUB 1-28)	80	62
49	17	10127	9942	gplI222606I H(+) - transporting ATP synthase [Streptomyces lividans] >pirSIS37545 H+ - transporting ATP synthase (EC 3.6.1.34) alpha chain - Streptomyces lividans	80	50
54	19	6858	6424	gplI235953I MIS1 gene product [Saccharomyces cerevisiae] >gplJ03724IYSCKMIS1A_1 S.cerevisiae mitochondrial C-1-Tetrahydrofolate synthase gene (MIS1). [Saccharomyces cerevisiae] >pirSIS28174 C1-tetrahydrofolate synthase precursor, mitochondrial - yeast (Saccharomyces	80	52
101	7	2602	3246	gplM96434I oxaloacetate decarboxylase [Salmonella typhimurium] >pirSISB44465 sodium ion pump oxaloacetate decarboxylase subunit alpha - Salmonella typhimurium	80	64
116	10	4191	3247	gplU32847I ribosomal protein S1 homolog, RNA-binding protein [Haemophilus influenzae]	80	63
170	3	2283	1651	gplU43739I FtsZ [Borrelia burgdorferi]	80	67
561	1	112	555	gplI45199I glucose-6-phosphate 1-dehydrogenase [Haemophilus influenzae] >gplU32737IHU32737_8 glucose-6-phosphate 1-dehydrogenase [Haemophilus influenzae] >gplU00074IHU00074_79 glucose-6-phosphate 1-dehydrogenase [Haemophilus influenzae] >gplU32846IHU32846_7 gluco	80	69

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

573	1	234	452	gpIU000191	B2235_C2_195 [Mycobacterium leprae]			80	65
5	10	6650	7417	gpID265621	plasmid copy number control protein pcnB [Escherichia coli] >pirSISS45212 plasmid copy number control protein - Escherichia coli			79	58
27	5	2427	2768	gpIZ497821	peptide chain release factor 1 [Bacillus subtilis] >pirSISS53437 peptide chain release factor 1 - Bacillus subtilis			79	61
42	17	9764	10186	gpIX530721	S typhimurium rplJ and rplL genes for ribosomal protein L10 and L7/L12. [Salmonella typhimurium] >pirSIR5EB12 ribosomal protein L7/L12 - Salmonella typhimurium			79	61
57	2	428	1900	gpIJ048361	M.barkeri ATPase alpha and beta subunit (atpA and atpB) genes, complete cds. [Methanosaerina barkeri] >pirSIA34283 H+-transporting ATP synthase (EC 3.6.1.34) alpha chain - Methanosaerina barkeri			79	63
60	6	4417	3140	gpIJ048361	M.barkeri ATPase alpha and beta subunit (atpA and atpB) genes, complete cds. [Methanosaerina barkeri] >pirSIA34283 H+-transporting ATP synthase (EC 3.6.1.34) alpha chain - Methanosaerina barkeri			79	64
85	4	1001	1450	gpIL454511	carbon storage regulator [Haemophilus influenzae] >gpiU32763 HTU32763_8 carbon storage regulator [Haemophilus influenzae] >gpiU00077 HTU00077_72 carbon storage regulator [Haemophilus influenzae] >gpiU32709 HTU32709_9 carbon storage regulator [Haemophilus			79	49
123	2	1158	550	gpIM802151	uvrA02 protein [Streptococcus pneumoniae] >pirSIA42385 uvr-402 protein - Streptococcus pneumoniae plasmid pSB470			79	62
165	5	1287	1766	gpIL452621	ATP-binding protein [Haemophilus influenzae] >gpiU32744 HTU32744_7 ATP-binding protein [Haemophilus influenzae] >gpiU00075 HTU00075_40 ATP-binding protein [Haemophilus influenzae] >gpiU32690 HTU32690_8 transport ATPase [Haemophilus influenzae] >pirSIC640			79	54
6	2	1617	430	gpIX158671	Micrococcus luteus homolog of the E.coli uvrA gene. [Micrococcus luteus] >pirSIS04781 uvrA protein - Micrococcus luteus			78	65
21	11	7812	9110	gpIL088541	valyl-tRNA synthetase [Lactobacillus casei] >pirSIS049856 valine--			78	58

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				tRNA ligase (EC 6.1.1.9) - Lactobacillus casei		
22	21	14655	14203	gpiU000391 [E. coli chromosomal region from 76.0 to 81.5 minutes. [Escherichia coli] >pirIS47728 yhiD protein - Escherichia coli >gplD11109]ECO10KLS_4 ORF-C [Escherichia coli] {SUB 1-148}	78	66
29	24	12039	10612	gpiD640061 hypothetical protein [Synechocystis sp.]	78	62
44	8	4761	4459	gpiI2542I phosphatase [Treponema denticola]	78	64
84	9	5019	4426	gpiU20445I BirA protein [Bacillus subtilis]	78	47
100	14	3613	3365	gpiU4507I glutamine amidotransferase [Haemophilus influenzae]	78	61
				>gpiU32726IHTU32726_10 glutamine amidotransferase [Haemophilus influenzae] >gpiU00073IHTU00073_51 glutamine amidotransferase [Haemophilus influenzae]		
				>gpiU32935IHTU32835_5 glutamine amidotransferase [Haemophilus influenzae]		
122	7	3432	4901	gpiM8858II transfer RNA-Leu synthetase [Bacillus subtilis] >pirIS41882	78	67
172	1	26	2777	gpiX53456I leucine-tRNA ligase (EC 6.1.1.4) - Bacillus subtilis	78	42
726	1	421	140	gpiU3420I plasma membrane Ca2+-pump (PMCA1b) [Sus scrofa] >pirIS13057 Ca2+-transporting ATPase (EC 3.6.1.38) - pig	78	42
10	19	10346	9219	gpiM37487I cardiac triadin isoform 3 [Oryctolagus cuniculus] >pirIS28485IAT28485_1	78	69
				protein D [Haemophilus influenzae] >pirIS43576 protein D precursor - Haemophilus influenzae	77	58
30	13	8251	9300	gpiM96343I MreB protein [Bacillus subtilis]	77	53
36	20	13764	16661	gpiU3047I DNA polymerase III helicase [Vibrio cholerae]	77	62
42	26	17880	18242	gpiM22622I biflexa acetylomithine deacetylase (argE; complete cds), and ribosomal subunit protein (rpsL; 5' end) genes, [Leptospira biflexa] >pirIS1A31840 RNA polymerase beta chain homolog - Leptospira biflexa (serotype patoc)	77	60
				Leptospira biflexa (serotype patoc) polyprotein 1b [lactate dehydrogenase-elevating virus]	77	66
70	4	3305	3072	gpiU15146I CapD [Staphylococcus aureus]	77	60
74	1	3	890	gpiU10927I orf304 gene product [Treponema pallidum]	77	47
101	10	3523	4605	pirISB3650 oxaloacetate decarboxylase (EC 4.1.1.3) beta chain - Klebsiella pneumoniae	77	60

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

115	8	3077	3955	gpiX916551	IcpA gene product [Bacillus subtilis] >gpiD176301BACGPR_4	77	54
121	19	10031	8298	gpiU013761	ORF80 protein [Bacillus subtilis] [SUB 1-327]	77	58
					ATP-binding protein [Escherichia coli] >gpiU013761U01376_4		
					ATP-binding protein [Escherichia coli] >gpiM831381ECOFTSHIA_2		
					ftsH gene product [Escherichia coli] [SUB 4-647]		
137	3	1603	92	gpiM220391	CTP synthetase [Bacillus subtilis] >gpiZ49782IBSDNA320D_10	77	54
					CTP synthetase [Bacillus subtilis] >pinSISYB5TP CTP synthetase		
					(EC 6.3.4.2) - Bacillus subtilis		
160	1	3	338	gpiZ121601	gidA gene product [Borrelia burgdorferi]	77	60
					>gpiZ121601BBGIDAG_1 division protein [Borrelia burgdorferi]		
					[SUB 529-593] >gpiX95669BBTHDFGD_2 gidA gene product		
					[Borrelia burgdorferi] [SUB 1-29] >gpiX95668BBGIDMOXR_1		
					gidA gene product [Borrelia burgdorferi]		
195	1	388	2	gpiM647301	DNA mismatch repair protein [Escherichia coli]	77	59
					>gpiU29579IECU29579_30 DNA mismatch repair protein		
					[Escherichia coli]		
8	10	5202	5498	gpiU151861	sua5 [Mycobacterium leprae]	76	58
12	3	1872	2031	gpiM172821	elastin [Homo sapiens]	76	61
12	10	5367	5642	gpiL456611	transketolase 1 [TK 1] [Haemophilus influenzae]	76	60
					>gpiU327831HTU32783_2 transketolase 1 (TK 1) [Haemophilus		
					influenzae] >gpiU000791HTU00079_82 transketolase 1 (TK 1)		
					[Haemophilus influenzae] >gpiU327291HTU32729_2 transketolase		
					2 [Haemophilus influenzae] >pi		
19	18	13976	13413	gpiL772461	YppQ gene product [Bacillus subtilis]	76	53
28	3	1978	1724	gpiX782061	Glucose-6-phosphate isomerase [Leishmania mexicana]	76	52
42	18	10347	11399	gpiL484881	RNA polymerase beta subunit [Borrelia burgdorferi]	76	62
42	28	18795	19211	gpiL452211	ribosomal protein S7 [Haemophilus influenzae]	76	61
					>gpiU327391HTU32739_9 ribosomal protein S7 [Haemophilus		
					influenzae] >gpiU000741HTU00074_100 ribosomal protein S7		
					[Haemophilus influenzae] >gpiU328481HTU32848_9 30S		
					ribosomal protein S7 [Haemophilus influenzae]		

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

45	14	11488	12753	gpiU283771	metK gene product [Escherichia coli]	76	64
57	9	6841	7014	gpiL006731	tyrosine aminotransferase [Trypanosoma cruzi]	76	53
83	3	870	1112	gpiX603431	glyceraldehyde 3-phosphate dehydrogenase [Hordeum vulgare] >priSIDEBHG glyceraldehyde-3-phosphate dehydrogenase (EC 1.2.1.12) - barley	76	35
89	3	1288	1584	gpiM301981	recQ gene product [Escherichia coli]	76	55
100	11	2971	2693	gpiU375231	glucosamine synthetase [Sphingomonas yanoikuyae] >gpiU375231SYU37523_3 glucosamine synthetase	76	65
101	11	4526	4879	gpiM964361	oxaloacetate decarboxylase [Klebsiella pneumoniae] >priSIA44464 oxaloacetate decarboxylase beta subunit (C terminus) - Klebsiella pneumoniae (fragment) >gpiM262901KPNOADB2_1 oadB gene product [Klebsiella pneumoniae] [SUB 4-141]	76	61
150	16	12992	13774	gpiU000211	Mycobacterium leprae cosmid L247. [Mycobacterium leprae]	76	59
198	1	12	356	gpiL132921	ATPase [Enterococcus hirae]>priSIA45995 Cu2+-transporting ATPase (EC 3.6.1.-) - Enterococcus hirae	76	47
6	3	3386	1311	gpiL448941	excinuclease ABC subunit A [Haemophilus influenzae] >gpiU32711IHTU32711_1 excinuclease ABC subunit A [Haemophilus influenzae]>gpiU00071IHTU00071_64 excinuclease ABC subunit A [Haemophilus influenzae] >gpiU32820IHTU32820_1 excinuclease ATPase subunit [Hae	75	59
8	32	18366	18902	gpiL479711	ribosomal protein S5 [Bacillus subtilis]	75	54
8	35	21574	21807	gpiM264141	ribosomal protein S11 [Bacillus subtilis] >gpiL47971IBACRPLP_20 ribosomal protein S11 [Bacillus subtilis]>gpiM13957IBACRPOA_2 B_subtilis DNA sequence of the rpsM-rpoA interval. [Bacillus subtilis]>priSIR3BSS1	75	60
15	5	2866	3408	gpiX159811	E. coli sbcC gene (ORF 45) for SbcC. [Escherichia coli] >priSBVECSC sbcC protein - Escherichia coli >gpiM64787ECOARAJ_1 sbcC gene product [Escherichia coli] [SUB 378-1048]	75	62
22	23	14859	15518	gpiL458931	periplasmic serine protease D0 and heat shock protein	75	48

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				[Haemophilus influenzae] >gplU32805 HTU32805_12 periplasmic serine protease Dd and heat shock protein [Haemophilus influenzae] >gplU00082 HTU00082_28 periplasmic serine protease Dd and heat shock prote		
24	4	1217	999	gplJ05510 Rat inositol-1,4,5-triphosphate receptor mRNA, complete cds. [Rattus norvegicus] >priSIB36579 inositol 1,4,5-triphosphate receptor 2 - rat >gplU38665 RN38665_1 inositol 1,4,5-triphosphate receptor [Rattus norvegicus] (SUB 1612-1859)	75	58
26	1	266	3	gplY00402 Drosophila melanogaster mRNA for phosphoenolpyruvate carboxykinase (GTP) (PEPCK, EC 4.1.1.32). [Drosophila melanogaster] >priS QYFFGM phosphoenolpyruvate carboxykinase (GTP) (EC 4.1.1.32) precursor, mitochondrial - fruit fly (Drosophila melanogaster)	75	43
28	11	5554	4961	gplL45197 putative glucose-6-phosphate dehydrogenase isozyme [Haemophilus influenzae] >gplU32737 HTU32737_6 putative glucose-6-phosphate dehydrogenase isozyme [Haemophilus influenzae] >gplU00074 HTU00074_77 putative glucose-6-phosphate dehydrogenase isozyme [Haemophilus influenzae]	75	60
28	13	6156	5662	gplL45199 glucose-6-phosphate 1-dehydrogenase [Haemophilus influenzae] >gplU32737 HTU32737_8 glucose-6-phosphate 1-dehydrogenase [Haemophilus influenzae] >gplU00074 HTU00074_79 glucose-6-phosphate 1-dehydrogenase [Haemophilus influenzae] >gplU32846 HTU32846_7 gluco	75	56
35	38	19740	21392	gplD84214 YbbQ [Bacillus subtilis]	75	51
37	1	799	2	gplL29053 HtpG gene product [Vibrio fischeri]	75	58
37	3	922	2097	gplX70943 aspartyl-tRNA synthetase [Thermus aquaticus thermophilus] >priS33743 aspartate-tRNA ligase (EC 6.1.1.12) - Thermus aquaticus	75	58
37	16	10484	9618	gplD64002 hypothetical protein [Synechocystis sp.]	75	56
38	2	348	701	gplU14345 6-phosphogluconate dehydrogenase [Salmonella enterica]	75	60
41	8	2153	2428	gplU211571 sarcolemmal associated protein-3 [Oryctolagus cuniculus]	75	65
47	8	4123	4566	gplL39821 trxA gene product [Chlamydia psittaci]	75	52

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

48	11	8831	8241	ep[Z32850] [Ricinus communis]	pyrophosphate-dependent phosphofructokinase beta subunit	75	56
54	26	8599	8018	epIM609171	purine nucleoside phosphorylase [Escherichia coli] >gplU14003[ECOJW93_295 purine-nucleoside phosphorylase [Escherichia coli]>pirSIA27854 purine-nucleoside phosphorylase (EC 2.4.2.1) - Escherichia coli	75	50
64	25	10848	11318	gpIX581141	testis-specific RNA [Drosophila hydei]	75	62
91	2	165	10251	gpIM885811	transfer RNA-Leu synthetase [Bacillus subtilis]>pirSIA41882	75	66
113	3	1967	26833	gpIX647051	leucine-tRNA ligase (EC 6.1.1.4) - Bacillus subtilis secA gene product [Antithamnion sp.]>pirSIS42707 secA protein	75	60
137	11	4611	5657	gpIL447111	-red alga (Antithamnion sp.)	75	61
					DNA mismatch repair protein [Haemophilus influenzae] >gplU32692[HTU32692_3 DNA mismatch repair protein [Haemophilus influenzae]>pirU00069[HTU00069_65 DNA mismatch repair protein [Haemophilus influenzae] >gplU32801[HTU32801_3 DNA mismatch repair protein [75	
228	1	20	397	ep[U03991]	beta-galactosidase alpha peptide [Cloning vector pUC1918] >gplU33186[CVU33186_1 beta-galactosidase alpha peptide [Cloning vector pSUM36] [SUB 41-96]	75	72
245	1	251	469	ep[Z542061]	BTF [Bovine herpesvirus 1]	75	
420	1	81	236	gpIL236511	C. elegans cosmid C29E4. [Caenorhabditis elegans] >pirSIS44767 C29E4.1 protein - Caenorhabditis elegans	75	75
458	1	184	2	gpIL061471	golgin-95 [Homo sapiens]>pirSJH0821 95K golgi antigen - human	75	67
8	24	13252	13551	gpIX549941	ribosomal protein S19 [Bacillus stearothermophilus] >pirSR3BS19 ribosomal protein S19 - Bacillus stearothermophilus	74	56
8	33	18841	19551	gpIL479711	ribosomal protein L15 [Bacillus subtilis] >gplD00619[BACSECY_2_B.subtilis sec Y gene. [Bacillus subtilis] >gplM31102[BACSPCR_3_B.subtilis spectinomycin resistance (spc) genes, complete cds. [Bacillus subtilis] >gplM31102[BACSPCR_3_Bacillus subtilis spectin	74	55

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

18	7	6261	5623] gpiX68709]	whiG-Stv gene product [Streptovorticillium griseocarneum]	74	55
21	16	10638	10162 gpiU189971	>pirSIS29615 whiG protein - Streptovorticillium griseocarneum Escherichia coli K-12 chromosomal region from 67.4 to 76.0 minutes. [Escherichia coli]	74	54
26	25	14169	14657 gpiU39691I	methylgalactoside permease ATP-binding protein [Mycoplasma genitalium] >pirSIB64213 methylgalactoside permease ATP-binding protein homolog - Mycoplasma genitalium (SGC3) >gpiU02149 MGU02149_1 Mycoplasma genitalium random genomic clone sc8a, partial cds.	74	52
29	16	6788	6084 gpiU00013I	ppsl [Mycobacterium leprae]	74	46
29	17	7660	6170 gpiU00013I	ppsl [Mycobacterium leprae]	74	54
41	5	1572	1955 gpiL15191II	phosphoenolpyruvate:sugar phosphotransferase system enzyme I [Streptococcus mutans] >gplL15191ISTRPHOSPHO_2 phosphoenolpyruvate:sugar phosphotransferase system enzyme I [Streptococcus mutans]	74	60
70	2	856	1119 gpiX79146I	ImbS gene product [Streptomyces lincolnensis] >pirSIS44965 ImbS protein - Streptomyces lincolnensis	74	56
79	2	893	1453 gpiL76303I	flagellar basal body rod protein [Borrelia burgdorferi]	74	50
80	1	325	2 gpiM80473I	uvr/dnaA gene product [Bacillus subtilis] >gpiM64048IBACDINA76_2 Bacillus subtilis DNase inhibitor (dinA76) gene, complete cds and promoter region. [Bacillus subtilis] [SUB 1-57]	74	58
82	4	2176	1022 gpiS56812I	phosphoenol pyruvate carboxykinase...F-ATPase epsilon subunit [Chlorobium limicola, Genomic, 4 genes, 5477 nt]. [Chlorobium limicola]	74	64
100	18	4634	4086 gpiL45070I	glutamine amidotransferase [Haemophilus influenzae] >gpiU32726 HTU32726_10 glutamine amidotransferase [Haemophilus influenzae] >gpiU00073 HTU00073_51 glutamine amidotransferase [Haemophilus influenzae] >gpiU32835 HTU32835_5 glutamine amidotransferase [Hae	74	58
101	6	1817	2752 gpiJ03885I	K.pneumoniae oxalacetate decarboxylase alpha subunit gene, complete cds. [Klebsiella pneumoniae] >pirSIA28088 oxaloacetate	74	59

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

101	19	9164	10120	gpiU000241	decarboxylase (EC 4.1.1.3) alpha chain - Klebsiella pneumoniae	74	53
117	1	913	305	gpiL188671	u000242 [Mycobacterium tuberculosis] ATP-dependent protease ATPase subunit [Escherichia coli] >priSIA48709 ATP-dependent Clp proteinase (EC 3.4.-.-)	74	52
141	6	3736	4149	gpiU189971	regulatory chain X - Escherichia coli Escherichia coli K-12 chromosomal region from 67.4 to 76.0	74	52
159	2	1170	1460	gpiX637571	spoIIA gene product [Bacillus megaterium] >priSIA48402 stage II spongulation protein spoIIAA - Bacillus megaterium	74	44
604	1	2	844	gpiU283771	metK gene product [Escherichia coli]	74	60
4	4	3687	3995	gpiU441181	ribosomal protein L20 [Pseudomonas syringae pv. syringae]	73	44
8	23	12415	13254	gpiZ216771	ribosomal protein L2 [Thermotoga maritima] >priSIS40191	73	63
11	25	12075	13619	gpiZ186311	ribosomal protein L2 - Thermotoga maritima ORF2 gene product [Bacillus subtilis] >priSIC36905 nusA	73	50
19	11	8644	9660	gpiM636551	homolog - Bacillus subtilis delta-2-isopentenyl pyrophosphate transferase [Escherichia coli]	73	48
					>gpiU14031ECOUW93_83 tRNA delta-2-isopentenylpyrophosphate (IPP) transferase [Escherichia coli] >priSIB37318 delta(2)-isopentenylpyrophosphate transferase (EC 2.5.1.-) - Escherichia coli >		
22	5	4833	3916	gpiU189971	Escherichia coli K-12 chromosomal region from 67.4 to 76.0	73	52
23	2	32	616	gpiX743571	minutes. [Escherichia coli] skp gene product [Pasteurella multocida] >priSIS47341 skp	73	42
26	4	2728	3381	gpiL450341	hypothetical protein (SP:P31216) [Haemophilus influenzae] >gpiU32723[H TU32723_4 hypothetical protein (SP:P31216) [Haemophilus influenzae] >gpiU00073[H TU00073_15 hypothetical protein (SP:P31216) [Haemophilus influenzae]]	73	57
					>gpiU32832[H TU32832_2 GTPase [Haemop 3'-exo-deoxyribonuclease [Bacillus subtilis]]		
26	7	4429	4824	gpiD261851	ORF246 gene product [Escherichia coli]	73	57
28	2	1598	846	gpiX595511	>gpiD101651ECORUVIC_3 Orf26 [Escherichia coli] >priSIC38113	73	52

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

30	10	4632	5738	gpIU1759[1]	26K hypothetical protein (nuvC 5' region) - Escherichia coli	73	53
34	3	1318	2271	gpIX820[72]	primary sigma factor [Borrelia burgdorferi] orf2 gene product [Pseudomonas aeruginosa] >pirSIS49379	73	53
					hypothetical protein 2 - Pseudomonas aeruginosa >pirSIS57900		
					hypothetical protein 2 - Pseudomonas aeruginosa		
38	1	2	571	gpIL451[94]	6-phosphogluconate dehydrogenase, decarboxylating [Haemophilus influenzae] >gpIU3273[1]HIU3273[7]_3'_6'- phosphogluconate dehydrogenase, decarboxylating [Haemophilus influenzae] >gpIU0074[HIU0074]_74'_6'-phosphogluconate dehydrogenase, decarboxylating [Haemophilus influenzae] >gpIU0074[HIU0074]_74'_6'-phosphogluconate	73	59
42	15	8010	9143	gpIZ1183[9]	dehydrogenase, decarboxylating [Haemophilus influenzae] >gpIU0074[HIU0074]_74'_6'-phosphogluconate	73	49
51	3	4067	4717	gpIL252[88]	ribosomal protein L1 [Thermotoga maritima] >pirSIR5HG1T	73	49
56	6	2887	3738	gpID641[16]	ribosomal protein L1 - Thermotoga maritima	73	54
57	4	2201	3697	gpIX795[16]	gyrase-like protein alpha subunit [Staphylococcus aureus] ORF4 [Bacillus subtilis]	73	50
					membrane ATPase [Halofera volcanii] >pirSIS4514S H+ transporting ATP synthase (EC 3.6.1.34) beta chain - Halofera volcanii >pirSIS55896 membrane ATPase B chain - Halofera volcanii		
57	7	5549	5352	gpIU0782[5]	histone variant H1.1(a) [Parechinus angulosus]	73	31
79	3	1482	1985	gpIU437[39]	FlgC [Borrelia burgdorferi] >gpIL7630[3]BORFTSA_8 flagellar basal body rod protein [Borrelia burgdorferi]	73	57
100	30	12237	11806	gpID6400[3]	hypothetical protein [Synechocystis sp.] >gpID6400[3]SYCSLLE_14 hypothetical protein [Synechocystis sp.]	73	50
142	11	13952	13527	gpIX756[27]	spoIII E gene product [Coxiella burnetii] >pirSIS4313Z spoIII E protein - Coxiella burnetii >pirSIS31759 hypothetical protein 274 - Coxiella burnetii [SUB 505-778]	73	35
167	4	1451	2062	gpIX864[8][1]	nrf gene product [Clostridium perfringens]	73	50
668	2	143	373	pirSIS3371[6]	proline dehydrogenase - Salmonella typhimurium	73	60
2	11	7849	9303	gpIL3184[5]	UDP-N-acetyl muramate-alanine ligase [Bacillus subtilis]	72	45
5	6	3790	4449	gpIM2589[9]	antimicrobial protein [Artificial gene]	72	38
6	7	5180	4983	gpIX7285[7]	traX gene product [Streptomyces coelicolor] >pirSIS39854	72	63

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				[hypothetical protein X - Streptomyces coelicolor] >pirIS32232		
8	29	16835	16990	gpIZ2282031	[hypothetical protein X - Streptomyces coelicolor]	
				TOR2 gene product [Saccharomyces cerevisiae] >pirIS38040	72 55	
				TOR2 protein - yeast (Saccharomyces cerevisiae)		
16	14	10815	9463	gpID640051	[hypothetical protein [Synedocystis sp.]	
18	5	5680	5177	gpIX687091	whiG-SIV gene product [Streptovorticillium griseocarneum]	
				>pirIS29615 whiG protein - Streptovorticillium griseocarneum	72 46	
25	4	1968	1804	epIM220391	fructose-bisphosphate aldolase [Bacillus subtilis]	
				>gpIZ497821BSDNA320D_13 fructose bisphosphate aldolase	72 51	
				[Bacillus subtilis] >pirSID32354 fructose-bisphosphate aldolase		
				(EC 4.1.2.13) - Bacillus subtilis >gpIS425901S42590_1 fructose		
				1,6-bisphosphate aldol		
26	27	15501	16058	gpIU396911	ribose transport system permease protein [Mycoplasma genitalium]	
				>pirIS64213 ribose transport system permease homolog -	72 48	
				Mycoplasma genitalium (SGC3)		
27	3	1873	956	gpIZ496331	S.cerevisiae chromosome X reading frame ORF YJR133w.	
				[Saccharomyces cerevisiae] >pirIS57156 hypothetical protein	72 48	
35	8	4530	5426	gpIM143391	YJR133w - yeast (Saccharomyces cerevisiae)	
				S.pneumoniae DpnII gene region encoding dpmM, dpmA, dpmB,	72 59	
				complete cds. [Streptococcus pneumoniae]		
				>gpIM112261STRDPNM_1 S.pneumoniae DpmM gene encoding		
				Dpn II DNA methylase, complete cds. [Streptococcus		
				pneumoniae] >gpIM143391STRDPN2A_2 dpmM gene product		
35	31	17563	18165	gpIZ111651	641 aa (68 kD) gene product of ORF641. [Rhodobacter capsulatus]	
				>pirISG28771 hypothetical protein C2814 (photosynthetic gene	72 57	
				cluster) - Rhodobacter capsulatus		
36	15	92228	10505	gpIU397031	glycyl-tRNA synthetase [Mycoplasma genitalium] >pirISG64227	
				glycine-tRNA ligase (EC 6.1.1.14) - Mycoplasma genitalium	72 51	
62	9	3780	4169	gpIM645191	(SGC3)	
				transport protein [Escherichia coli] >pirISIA40840	72 62	
69	11	6417	5293	epIX040271	spermidine/putrescine transport protein A - Escherichia coli	
				E. coli mutD/dnaQ-mh region for DNA polymerase III epsilon	72 63	

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

					subunit and RNAase H. [Escherichia coli] >gpiK005521ECORNH_1 E.coli mh gene coding for ribonuclease H. [Escherichia coli] >gpiK009851ECORNHQ_2 ribonuclease H [Escherichia coli] >gpiV003371ECRNH	
79	12	6913	7581	gpIM727181	B.subtilis ftaA locus operon. [Bacillus subtilis] >gpIX56049IBSFLAAO_3 B. subtilis FtaA locus operon. [Bacillus subtilis] >pirSIIPWB5AS H+-transporting ATP synthase alpha chain homolog - Bacillus subtilis	72 53
82	2	475	281	gpIZ150251	Bat2 gene product [Homo sapiens] >pirSISS37671 bat2 protein - human Drosophila melanogaster mRNA for phosphoenolpyruvate carboxykinase (GTP) (PEPCK, EC 4.1.1.32). [Drosophila melanogaster] >pirSIQYFFGM phosphoenolpyruvate carboxykinase (GTP) (EC 4.1.1.32) precursor, mitochondrial - fruit fly [Drosophila melanogaster]	72 72
82	3	1136	480	gpIY004021	Drosophila melanogaster PcrA - Staphylococcus aureus >pirSISS27667 DNA helicase pcrA - Staphylococcus aureus >pirSISS39923 pcrA protein - Staphylococcus aureus	72 57
89	1	1041	4	gpIM6311761	PPI-dependent phosphofructo-1-kinase [Entamoeba histolytica] >gpIU12513IEHU12513_1 PPI-dependent phosphofructo-1-kinase [Entamoeba histolytica] >pirSISS52082 PPI-dependent phosphofructo-1-kinase - Entamoeba histolytica	72 61
101	4	1838	1557	gpID252201	selenoprotein P like protein [Bos taurus] >gpID252201BOVSPP_1 selenoprotein P like protein [Bos taurus]	72 44
105	4	2929	1649	pirSI SYBS YF	tyrosine-tRNA ligase (EC 6.1.1.) - Bacillus stearothermophilus	72 54
116	5	930	1568	gpIX819901	leader peptidase I [Phormidium laminosum] >pirSISS1921 leader peptidase (EC 3.4.99.36) - oscillatoriacean cyanobacterium	72 50
134	2	1073	477	gpIU332101	methyl accepting chemotaxis homolog [Treponema denticola]	72 54
136	1	2	268	gpIU151401	ribosomal protein IF-1 [Mycobacterium bovis]	72 52
144	2	119	310	gpIL288101	regulatory protein [Aspergillus nidulans]	72 61
161	2	787	464	gpIM58002	bacterial cell wall hydrolase [Streptococcus faecalis] >pirSI38109	72 39

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				autolysin - Enterococcus faecalis		
186	2	229	603	gplID640001 hypothetical protein [Synechocystis sp.]	72	59
555	2	157	456	gplX775151 pyruvate oxidoreductase [Rhodospirillum rubrum] >pirsIS41961	72	49
5	2	858	1499	gpm113301 pyruvate oxidoreductase - Rhodospirillum rubrum	71	52
				>gplD835361ECOTSF_5 CDP-diglyceride synthetase [Escherichia coli]	71	
				>pirsISYECDF phosphatidate cytidylyltransferase [Escherichia coli]	71	
11	28	14189	16321	gpxX005131 CDP-diglyceride synthetase [Escherichia coli] >pirsISYECDF phosphatidate cytidylyltransferase (EC 2.7.7.41) - Escherichia coli	71	54
				E.coli nusA operon including genes for Met-tRNA-F2 (metY), 15	71	
				kd protein, NusA Protein (nusA), and initiation factor IF2 (infB). [Escherichia coli] >gpiU18997IECOUW67_98 protein chain initiation factor 2 [Escherichia coli]	71	
13	8	5270	3489	gplL459341 hypothetical protein (GB:U14003_302) [Haemophilus influenzae] >gpiU32809 HTU32809_12 hypothetical protein (GB:U14003_302) [Haemophilus influenzae]	71	50
				>gpiU00082 HTU00082_66 hypothetical protein (GB:U14003_302) [Haemophilus influenzae]	71	
16	35	21358	20309	gpm863511 triacylglycerol acylhydrolase [Streptomyces sp.] >pirsIJN0490	71	53
				28K lipase precursor - Streptomyces sp. (strain M11)	71	
21	7	5991	4771	gplD164371 PacS [Synechococcus sp.] >pirsIS36/41 cation-transporting ATPase pacS - Synechococcus sp.	71	49
29	11	4929	4390	gplID640041 hypothetical protein [Synechocystis sp.] >gplD640041SYCSLRF_7 hypothetical protein [Synechocystis sp.]	71	43
35	37	18671	19465	gplL460711 hypothetical protein (SP:P26242) [Haemophilus influenzae] >gpiU32822 HTU32822_12 hypothetical protein (SP:P26242) [Haemophilus influenzae] >gpiU00084 HTU00084_34 hypothetical protein (SP:P26242) [Haemophilus influenzae]	71	56
				>gpiU32768 HTU32768_15 transketolase D9651_10 gene product [Saccharomyces cerevisiae]	71	41
41	13	5053	5550	gpxX059911 Drosophila Cs gene, [Drosophila melanogaster] >pirsIS01103	71	64

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

42	11	5125	5469	gp L46067	hypothetical protein 2 - fruit fly (<i>Drosophila melanogaster</i>) H. influenzae predicted coding region H1435 [Haemophilus influenzae] >gp U32822 HTU32822_8 H. influenzae predicted coding region H1435 [Haemophilus influenzae] >gp U00084 HTU00084_30 H. influenzae predicted coding region H1435 [Haemophilus influenzae]	71	50
42	24	16051	17931	gp U000061	DNA-directed RNA polymerase, beta'-subunit [Escherichia coli] >gp X04642 STRPOB_2 S. typhimurium rpoB gene for RNA polymerase beta subunit. [Salmonella typhimurium] [SUB 1-20]	71	53
47	21	12695	12396	gp X74933	tRNA(mG3)methyltransferase [Salmonella typhimurium] >pir SIS37175 tRNA (guanine-N1->methyltransferase [EC 2.1.1.31]) - Salmonella typhimurium	71	40
48	4	2455	2279	gp Z22915	T26G10.1 [Caenorhabditis elegans]>pir SIS40731 hypothetical protein - Caenorhabditis elegans	71	52
48	7	5241	4102	gp M91598	3-phosphoglycerate kinase [Yarrowia lipolytica]	71	46
52	5	6567	4714	gp X73141	hemolysin [Serpulina hydysenteriae]	71	48
57	6	6423	4474	gp L45709	hypothetical protein (SP:P37024) [Haemophilus influenzae] >gp U32787 HTU32787_7 hypothetical protein (SP:P37024) [Haemophilus influenzae]>gp U00080 HTU00080_36 hypothetical protein (SP:P37024) [Haemophilus influenzae]>pir SISG64165 hypothetical protein H	71	54
61	8	3121	3765	gp L45095	hypothetical protein (GB:D26185_102) [Haemophilus influenzae] >gp U32728 HTU32728_13 hypothetical protein (GB:D26185_102) [Haemophilus influenzae] >gp U00073 HTU00073_76 hypothetical protein (GB:D26185_102) [Haemophilus influenzae] >gp U32837 HTU32837_6 H	71	44
63	19	7365	7027	gp L45782	hypothetical protein (SP:P33995) [Haemophilus influenzae] >gp U32794 HTU32794_8 hypothetical protein (SP:P33995) [Haemophilus influenzae]>gp U00081 HTU00081_13 hypothetical protein (SP:P33995) [Haemophilus influenzae] >gp U32740 HTU32740_8 ATP/GTP-utili	71	53
63	21	9982	8696	gp M19488	S.typhimurium prsA gene encoding phosphoribosylpyrophosphate	71	49

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				synthetase, complete cds. [Salmonella typhimurium] >pirSKIEBRT ribose-phosphate pyrophosphokinase (EC 2.7.6.1) - Salmonella typhimurium	
65	5	1824	2435	gpiID55650 adenylylate cyclase [Anabaena cylindrica]	71 50
118	1	650	1324	gpiZ36878 putative nicotinate phosphoryltransferase [Saccharomyces cerevisiae] >pirSIS51845 probable nicotinate phosphoribosyltransferase (EC 2.4.2.11) - yeast (Saccharomyces cerevisiae) >gplL11274 YSCNPTTRPBX_1 nicotinate phosphoribosyltransferase [Saccharomyces cerevisiae]	71 59
131	2	1615	746	gpiM30942 transfer RNA-Ile synthetase [Tetrahymena thermophila] >pirSA42399 isoleucyl-tRNA synthetase, ileRS - Tetrahymena thermophila (SGC5)	71 52
150	9	8131	7304	gpiID26185 regulation of Spo01 and Ort283 (probable) [Bacillus subtilis] >gpiX62539 BSORIGS_10 B. subtilis genes rpmH, rpmA, 50kd, gidaA and gidB - [Bacillus subtilis] >pirSIS18080 hypothetical protein 5 - Bacillus subtilis	71 55
157	3	1725	1366	gpiL44678 hypothetical protein (SP:P05848) [Haemophilus influenzae] >gpiU32688 HTU32688_13 hypothetical protein (SP:P05848) [Haemophilus influenzae] >gpiU00069 HTU00069_32 hypothetical protein (SP:P05848) [Haemophilus influenzae] >gpiU32798 HTU32798_1 H. influenzae	71 50
2	13	13010	10815	gpiID26185 transcription-repair coupling factor [Bacillus subtilis]	70 50
4	3	2796	3344	gpiU32757 initiation factor 3 [Haemophilus influenzae] >gpiL45952 HEAH11318_1 initiation factor 3 [Haemophilus influenzae] (SUB 49-183)	70 46
8	21	10812	11444	gpiX67014 ribosomal protein L3 [Bacillus stearothermophilus] >pirSIS24363	70 46
16	8	5949	5515	gpiX73141 hemolysin [Serpulina hyodysenteriae]	70 36
23	7	3492	5213	gpiZ38002 serine hydroxymethyltransferase [Bacillus subtilis] >pirSIS49363	70 49
29	28	13092	13259	gpiX05790 alpha-D-galactosidase [Homo sapiens] >gpiX14448 HSGLA_1	70 50
				alpha-D-galactosidase A [Homo sapiens] >pirSISGBHUA alpha-	

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				galactosidase (EC 3.2.1.22) A precursor - human >gplM13571 HUMAGALAA_1 GLA gene product [Homo sapiens] [SUB 27.429]		
35	1	306	4	gpiU23514 F48E8.6 gene product [Caenorhabditis elegans]	70	50
35	6	3610	3909	gpiMT6547 acyl carrier protein [Chloroplast Cryptomonas phai] >pirsIC41609 acyl carrier protein acpA - Cryptomonas sp. chloroplast (strain Phi) [SUB 2-81]	70	45
35	36	18773	18552	gpiX8733 ORF OR26.23 gene product [Saccharomyces cerevisiae] >gplS69545 S69545_1 Dhs1 [Saccharomyces cerevisiae] [SUB 220-702]	70	35
37	6	1980	2591	gpiX70943 aspartyl-tRNA synthetase [Thermus aquaticus thermophilus] >pirsS33743 aspartate-tRNA ligase (EC 6.1.1.12) - Thermus aquaticus	70	53
39	8	3707	3195	gpiU32699 acid phosphatase (?) [Haemophilus influenzae]	70	47
44	5	1398	2432	gpiX73140 hemolysin [Serpulina hyodysenteriae]	70	49
52	3	4080	2803	gpiU10400 YHR011w gene product [Saccharomyces cerevisiae] >pirsS46786 serine-tRNA ligase homolog - yeast (Saccharomyces cerevisiae)	70	46
56	9	3984	5726	gpiD64116 ORF4 [Bacillus substillis]	70	48
62	8	3523	4071	gpiL45980 spermidine/putrescine transport ATP-binding protein [Haemophilus influenzae] >pilU32813 HTU32813_12 spermidine/putrescine transport ATP-binding protein [Haemophilus influenzae] >pilU00083 HTU00083_23 spermidine/putrescine transport ATP-binding protein [Haemophilus influenzae] >pilU32813 HTU32813_12	70	48
74	5	1229	1444	gpiU41536 F56E3_3 gene product [Caenorhabditis elegans]	70	35
77	12	6229	5666	gpiX03038 E. coli adk gene for adenylylate kinase. [Escherichia coli] >gplM38777 ECOAPTADK_6 E.coli sequence of the apt-adk region. [Escherichia coli] >pirsIKIECA adenylylate kinase (EC 2.7.4.3) - Escherichia coli >gplD90259 ECOADKVIS_1 E.coli adk, visA genes and ORFs	70	47
81	14	13662	13255	pirsIR3BS9 ribosomal protein S9 - Bacillus stearothermophilus	70	54
100	12	3525	2704	gpiL45070 glutamine amidotransferase [Haemophilus influenzae]	70	54

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				>gpiU32726 HTU32726_10 glutamine amidotransferase [Haemophilus influenzae] >gpiU00073 HTU00073_5I glutamine amidotransferase [Haemophilus influenzae]		
100	19	5203	4832	>gpiU14003 Escherichia coli K-12 chromosomal region from 92.8 to 00.1 minutes. [Escherichia coli] >pirlSIS56432 hypothetical protein o259a - Escherichia coli >pirlSIS46294 peptidylprolyl isomerase (EC 5.2.1.8) - Escherichia coli [SUB 55-259]	70	56
100	29	11866	11027	>gpiD64003 hypothetical protein [Synechocystis sp.] sp.]	70	56
101	15	7124	6945	>gpiX63765 ribosomal protein RL9 [Synechococcus sp.] >pirlSIS22206 ribosomal protein L9 - Synechococcus sp.	70	50
121	22	10796	12685	>gpiX94607 Mef1 protein [Saccharomyces cerevisiae]	70	56
123	3	2235	1128	>gpiM80215 lvsA02 protein [Streptococcus pneumoniae] >pirlSIA42385 uvr- 402 protein - Streptococcus pneumoniae plasmid pSB470	70	51
143	3	3370	872	>gpiL45958 ion protease [Haemophilus influenzae] >gpiU32812 HTU32812_1 ion protease [Haemophilus influenzae] >gpiU00083 HTU00083_2 ion protease [Haemophilus influenzae] >gpiU32757 HTU32757_9 Lon/Sms-related endopeptidase (no ATPase domain) [Haemophilus influenzae] >	70	34
175	2	852	118	>gpiL45009 protein E [Haemophilus influenzae] >gpiU32721 HTU32721_3 protein E [Haemophilus influenzae] >gpiU00072 HTU00072_93 protein E [Haemophilus influenzae] >gpiU32830 HTU32830_3 essential protein [Haemophilus influenzae] >pirlSH4063 protein E (gpcE) homolog -	70	50
313	1	474	151	>gpiD10388 N-acetyl muramoy-L-alanine amidase [Bacillus subtilis] >gplM81324 BACWL_B 1 N-acetyl muramoy-L-alanine amidase [Bacillus subtilis] >gplM8764 BACLYTABCD_4 amidase [Bacillus subtilis] >pirlSIB41322 N-acetyl muramoy-L-alanine amidase (EC 3.5.1.28) 50K precu	70	44
441	1	2	232	>gpiU14333 Amt [Mus musculus]	70	47
590	1	3	338	>gpiM2773 D. discoideum Thy 1 gene, complete cds. [Dictyostelium	70	56

TABLE 2. *Treponema pallidum* - Putative coding regions of novel proteins similar to known proteins

				discoideum] >pirSIS1YXD0TC thymidylate synthase-complementing protein - slime mold [Dictyostelium discoideum]
641	1	51	575	gpIL348791 NiTS gene product [Anabaena azollae] >gpIL348791ANA AZNIF_3
4	13	8224	9780	gpIX545481 NiTS gene product [Anabaena azollae] serine protease [Salmonella typhimurium] >pirSIS15337 heat shock protein htrA - Salmonella typhimurium
13	6	2524	3492	gpIL460001 threonyl-tRNA synthetase [Haemophilus influenzae] >gpIU32816HTU32816_6 threonyl-tRNA synthetase [Haemophilus influenzae] >gpIU00831HTU0083_43 threonyl-tRNA synthetase [Haemophilus influenzae] >gpIU32762HTU32762_4 threonyl-tRNA synthetase [Haemophilus
16	13	9323	7899	gpID640051 hypothetical protein [Synechocystis sp.]
16	50	29876	29175	gpIX520941 flagG protein product (AA 1-260) [Salmonella typhimurium] >pirSIXMEBFG basal body rod protein flagG - Salmonella typhimurium
16	51	30892	29936	gpIX520941 flagG protein product (AA 1-260) [Salmonella typhimurium] >pirSIXMEBFG basal body rod protein flagG - Salmonella typhimurium
16	56	34171	33398	gpID640061 hypothetical protein [Synechocystis sp.]
17	6	3294	2881	gpIL092281 Bacillus subtilis spoVA to serA region. [Bacillus subtilis] >pirSIS45550 hypothetical protein X8 - Bacillus subtilis
28	8	4806	36835	gpIL173201 acetate kinase [Bacillus subtilis] >pirSIB49935 acetate kinase (EC 2.7.2.1) - Bacillus subtilis
29	7	3371	2706	gpID640021 hypothetical protein [Synechocystis sp.]
29	19	8526	76033	gpID640041 hypothetical protein [Synechocystis sp.] >gpID640041SYCSLRF_5 hypothetical protein [Synechocystis sp.]
30	5	2460	2894	gpIU131651 DNA primase [Listeria monocytogenes]
44	2	130	1539	gpIX73140 hemolysin [Serpulina hydysenteriae]
44	9	5167	4751	gpIL25421 phosphatase [Treponema denticola]
45	7	7818	8294	gpIV03161 tRNA synthetase [Saccharomyces cerevisiae]

TABLE 2. *Treponema pallidum* - Putative coding regions of novel proteins similar to known proteins

				>gpJU01339 YSCMESL_1 Yeast (<i>S.cerevisiae</i>) methionyl-tRNA synthetase (mes1) gene, complete cds. [<i>Saccharomyces cerevisiae</i>]		
				>gpJU01339 YSCMESL_1 <i>S.cerevisiae</i> methionyl-tRNA synthetase (mes1) gene, complete cds. [<i>Sa</i>]		
47	2	1724	822	gpJU19610 NADH oxidase [Serpulina hydrotsentierae]	69	48
60	4	2316	1218	gpJU04836 M.barkeri A TPase alpha and beta subunit [atpA and atpB] genes, complete cds. [Methanosc礼ina barkeri] >pirISB34283 H+-transporting ATP synthase (EC 3.6.1.34) beta chain - Methanosc礼ina barkeri	69	51
67	1	1916	180	pirISI0941 spoTIE protein - <i>Bacillus subtilis</i> >gpIM17445 BACSPIII_A_3 B.subtilis sporulation protein spoTIEA and spoTIEB genes, complete cds and open reading frame X, 3' end. [<i>Bacillus subtilis</i>] (SUB 536-787)	69	49
70	7	7470	5644	gpIX522227 E.coli fhaA gene for the transcriptional activator of the formate hydrogenlyase. [<i>Escherichia coli</i>] >gpIU29579 ECU29579_28 transcriptional activator of the formate hydrogenlyase system [Escherichia coli] >pirISI12079 transcriptional activator fhaA - Esch	69	52
79	4	1967	2356	gpIL76303 flagellar basal body rod protein [Borrelia burgdorferi]	69	36
80	2	356	1201	gpIX73124 ipa-52r gene product [<i>Bacillus subtilis</i>] >pirIS339707 hypothetical protein - <i>Bacillus subtilis</i>	69	48
81	20	17491	17204	gpIM63176 helicase [Staphylococcus aureus] >pirIS277667 DNA helicase pcra - Staphylococcus aureus >pirISI339923 pcra protein - Staphylococcus aureus	69	48
86	2	345	809	gpIZ221970 34CP [Chloroplast Arabidopsis thaliana] >pirISI36637 signal recognition particle 34CP protein precursor - <i>Arabidopsis thaliana</i>	69	59
86	3	616	1602	gpIX01818 E. coli tmrD operon and nearby regions. [<i>Escherichia coli</i>] >pirIS07178 hypothetical protein, 48K (rpsP 5' region) - <i>Escherichia coli</i>	69	46
98	3	3371	1197	gpU29668 polynucleotide phosphorylase [<i>Bacillus subtilis</i>]	69	52
100	15	4224	3529	gpIL10328 glutamine amidotransferase [<i>Escherichia coli</i>]	69	51
100	28	11096	10338	gpID64003 hypothetical protein [Synechocystis sp.]	69	53

TABLE 2. Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				>gp D64003 SYCSLLE_14 hypothetical protein [Synechocystis sp.]	
103	5	4051	4554	gp L13078I antizyme [Escherichia coli] >pirlSIA4829I ornithine decarboxylase inhibitor - Escherichia coli	69 53
140	2	165	599	gp X54548I serine protease [Salmonella typhimurium] >pirlSIS15337 heat shock protein htrA - Salmonella typhimurium	69 53
165	3	584	309	gp U14003I Escherichia coli K-12 chromosomal region from 92.8 to 00.1 minutes [Escherichia coli] >pirlSIS56374 hypothetical protein f342 - Escherichia coli	69 57
165	7	1946	2371	gp D83536I hypothetical 23.3 kd protein [Escherichia coli] >gp D15061 ECORRNHK12_5 ORF217 [Escherichia coli] [SUB 5-217]	69 44
376	2	224	628	gp U09229I transcription factor [Rattus norvegicus] >pirlSIB53689 homeotic protein CDP2 - rat (fragment)	69 58
3555	3	410	757	gp D64005I hypothetical protein [Synechocystis sp.]	69 62
629	1	3	692	gp M91593I Mycoplasma mycoides SRPM54 gene, complete cds. [Mycoplasma mycoides] >pirlSIS35480 hypothetical protein 1 - Mycoplasma mycoides (SGC3) >pirlSIS27590 hypothetical protein 1 - Mycoplasma mycoides (SGC3) (fragment) [SUB 2-422]	69 48
1	6	1876	2325	gp U22817I RuvB Protein [Thermus aquaticus thermophilus]	68 54
8	27	1595	16290	gp M81748I ribosomal protein L24 [Bacillus subtilis] >gp L47971 BACRPLP_5 ribosomal protein L24 [Bacillus subtilis] >gp X15664 BSSPC_4.B.subtilis S10/spc operon rpmC, rpsQ, rplN, rplX, rplE, rpsN genes. [Bacillus subtilis] >pirlSR5BS2B ribosomal protein L24 - Bacil	68 50
8	36	21789	22193	gp U30821I alpha subunit of RNA polymerase [Cyanelle Cyanophora paradoxa]	68 57
10	7	5142	4447	gp D64003I hypothetical protein [Synechocystis sp.] >gp D64003 SYCSLLE_79 hypothetical protein [Synechocystis sp.]	68 43
11	8	4625	5395	gp U39703I DNA helicase II [Mycoplasma genitalium] >pirlSII64226 DNA helicase II (mutB1) homolog - Mycoplasma genitalium (SGC3)	68 46

TABLE 2.
Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				>gplX61517 MGMG08_1 M.genitalium random genomic sequence MG08. [Mycoplasma genitalium] [SUB 277-345]	
12	7	4202	4507	gplU29580 Escherichia coli K-12 genome; approximately 62 minute region. [Escherichia coli]	68 34
16	28	16838	16275	gplX65028 rotamase [Synechococcus sp. (PCC 7942)] >pirlSICSYC42 peptide/prolyl isomerase (EC 5.2.1.8) - Synechococcus sp. (PCC 7942)	68 68
21	8	7156	5951	gplD164371 PacS [Synechococcus sp.] >pirlSIS36741 cation-transporting ATPase pacS - Synechococcus sp.	68 51
22	20	13872	13486	gplL463191 endonuclease III [Haemophilus influenzae] >gplU32842 HTU32842_1 endonuclease III [Haemophilus influenzae] >gplU00086 HTU00086_51 endonuclease III [Haemophilus influenzae] >gplU32788 HTU32788_7 endonuclease III [Haemophilus influenzae] >pirlSIG64136 endonu	68 53
24	2	479	1228	gplU09189 loricrin [Mus musculus] >gplM34398 MUSLRCNA_1 loricrin [Mus musculus] >pirlSA35628 loricrin - mouse	68 64
28	10	4950	4804	gplU02025 insulin-like growth factor binding protein 5 [Mus musculus] >pirlSA54259 insulin-like growth factor binding protein 5 - mouse (fragment) [SUB 1-111]	68 45
30	4	1681	2442	gplL45098 amino deoxychorismate lyase [Haemophilus influenzae] >gplU32728 HTU32728_16 amino deoxychorismate lyase [Haemophilus influenzae] >gplU00073 HTU00073_79 amino deoxychorismate lyase [Haemophilus influenzae] >gplU32837 HTU32837_9 H influenzae predicted coding	68 52
36	17	10774	11841	gplU12735 amino alcohol phosphotransferase [Glycine max]	68 54
48	10	8283	7114	gplZ32850 pyrophosphate-dependent phosphofructokinase beta subunit [Ricinus communis]	68 52
49	13	8601	9068	gplX64324 calquestrin [Rana esculenta] >pirlSIS22418 calquestrin - edible frog	68 47
61	11	5514	6017	gplL42534 pIeD gene product [Caulobacter crescentus] >gplL42554 CCRPLED_1 pIeD gene product [Caulobacter crescentus]	68 46

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

66	2	1223	720	gpiM/28883	Rabbit macrophage cationic peptide I (MCP-1) mRNA, complete cds. [Oryctolagus cuniculus] >gpiM/28072RABMCP1AA_1	68	45
					Rabbit macrophage cationic peptide I (MCP-1) gene, complete cds. [Oryctolagus cuniculus] >pirSIA45811 macrophage cationic peptide I precursor		
70	28	17845	18918	gplL231471	phosphotransacetylase [Methanoscarcina thermophila] >pirSIA49338 phosphate acetyltransferase (EC 2.3.1.8) - Methanoscarcina thermophila >gpiU501891MTU50189_1 phosphate acetyltransferase [Methanoscarcina thermophila] [SUB 1-22]	68	52
75	11	4924	4022	gpiM/269341	E.coli ansA-ORF1 gene pair, complete cds. [Escherichia coli] >pirSIIQQECA5 hypothetical 23K protein (ansA 3' region) - Escherichia coli	68	50
81	15	14103	13660	gplL460751	ribosomal protein L13 [Haemophilus influenzae] >gpiU32823IH TU32823_2 ribosomal protein L13 [Haemophilus influenzae] >gpiU00084IH TU00084_38 ribosomal protein L13 [Haemophilus influenzae] >gpiU32769IH TU32769_2 ribosomal protein L13 [Haemophilus influenzae]	68	48
104	1	99	389	gplD261851	unknown [Bacillus subtilis]	68	31
121	3	1239	772	gplL452631	polypeptide deformylase (formylmethionine deformylase) [Haemophilus influenzae] >gpiU32745IH TU32745_1 polypeptide deformylase (formylmethionine deformylase) [Haemophilus influenzae] >gpiU00075IH TU00075_42 polypeptide deformylase (formylmethionine deformyl)	68	52
135	1	128	478	gplL448841	protein-export membrane protein [Haemophilus influenzae] >gpiU32710IH TU32710_5 protein-export membrane protein [Haemophilus influenzae] >gpiU00071IH TU00071_55 protein-export membrane protein [Haemophilus influenzae]	68	45
136	2	1040	387	gplD640041	>gpiU32819IH TU32819_4 protein-export me hypothetical protein [Synechocystis sp.] >gplD640041SYCSLRF_28 hypothetical protein [Synechocystis sp.] recombinational DNA repair protein [Haemophilus influenzae]	68	59
142	6	8745	9422	gplL450841	>gpiU32728IH TU32728_2 recombinational DNA repair protein	68	47

TABLE 2.
Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

			[Haemophilus influenzae] >gpiU00073[H TU00073_65 recombinational DNA repair protein [Haemophilus influenzae]	
150	14	12885	11392 gplL37094l >gpiU32836[H TU32836_4 recombinase [Bartonella bacilliformis]	68 48
170	2	1730	1008 gplU43739l FtsZ [Borrelia burgdorferi]	68 44
186	1	2	2635 gplL44885l protein-export membrane protein [Haemophilus influenzae] >gpiU32710[H TU32710_6 protein-export membrane protein [Haemophilus influenzae] >gpiU00071[H TU00071_56 protein-export membrane protein [Haemophilus influenzae] >gpiU32819[H TU32819_5 protein-export me	68 36
323	1	184	468 gplV000328l recA gene product [Escherichia coli] >pirsSIRQECA recA protein - Escherichia coli (SUB 2-353)	68 47
5	7	4447	6303 gplS63246l CHL15 gene product [Saccharomyces cerevisiae]	67 32
13	4	1692	2564 gplV00291l E.coli thrS, infC, rplT, phoS, pheT and himA genes encoding threonyl-tRNA synthetase, initiation factor IF3, ribosomal protein L20, phenylalanyl-tRNA synthetase and the alpha-subunit of the host integration factor. [Escherichia coli] >pirsSYECCR1 threonine protein-glutamate methyl esterase (EC 3.1.1.61) - Salmonella typhimurium >pirsIA26119 protein-glutamate O-methyltransferase (EC 2.1.1.80) - Salmonella typhimurium (fragment) (SUB 277-306)	67 47
19	21	16547	15312 pirsIXYYEB ET mutS gene product [Azotobacter vinelandii] >pirsIA53296 DNA mismatch repair protein MutS - Azotobacter vinelandii	67 50
23	3	624	1487 gplM63007l cell division protein [Mycoplasma genitalium] >pirsIE64250 cell division protein ftsH - Mycoplasma genitalium (SGC3)	67 50
26	28	16025	16981 gplZ33126l membrane forming protein [Mycoplasma capricolum] >pirsS48611 hypothetical protein - Mycoplasma capricolum (SGC3) (fragment) (SUB 1-101)	67 32
28	14	7151	6051 gplL45169l glucose-6-phosphate 1-dehydrogenase [Haemophilus influenzae] >gpiU32737[H TU32737_8 glucose-6-phosphate 1-dehydrogenase [Haemophilus influenzae] >gpiU00074[H TU00074_79 glucose-6-phosphate 1-dehydrogenase [Haemophilus influenzae]	67 50

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

30	6	2930	3220	gpiM10040	B.subtilis dnaE gene encoding DNA primase, complete cds. [Bacillus subtilis] >gpiX03897 BSSIG43_2 Bacillus subtilis sigma 43 operon with P23-dnaE-rpoD genes (dnaE for DNA primase, rpoD for RNA polymerase). [Bacillus subtilis] >pirSIA22282	67	42
35	15	8921	9520	gpJ03762	E.coli thioredoxin reductase gene, complete cds. [Escherichia coli] >pirSIRDECT thioredoxin reductase (NADPH) (EC 1.6.4.5) - Escherichia coli >gpL21749 ECOCYDD_1 thioredoxin reductase [Escherichia coli] [SUB 244-321]	67	52
36	3	2166	3407	gpL45171	sialic acid coprotease [Haemophilus influenzae] >gpiU32735 HUU32735_5 sialic acid coprotease [Haemophilus influenzae] >gpiU00074 HUU00074_51 sialic acid coprotease [Haemophilus influenzae] >gpiU32844 HUU32844_5 sialic acid coprotease [Haemophilus influenzae] >pirSIIH6407	67	53
36	14	7522	9255	gpiM27221	glutamyl-tRNA synthetase [Rhizobium meliloti] >pirSISYRZET glutamate-tRNA ligase (EC 6.1.1.17) - Rhizobium meliloti	67	48
45	6	6922	7989	gpIV013161	tRNA synthetase [Saccharomyces cerevisiae] >gpiI01339 YSCMES1_1 Yeast (S.cerevisiae) methionyl-tRNA synthetase (mes1) gene, complete cds. [Saccharomyces cerevisiae] >gpiI01339 YSCMES1_1 S.cerevisiae methionyl-tRNA synthetase (mes1) gene, complete cds. [Sa	67	48
47	18	10834	10559	gpIX74121	flagellar biosynthetic protein [Bacillus subtilis] >pirSIS34714 flagellar protein fliB - Bacillus subtilis	67	48
51	4	4609	5790	gpIX942241	DNA gyrase subunit A [Mycobacterium smegmatis] >gpiX94224 MSGYRBA_4 DNA gyrase subunit A [Mycobacterium smegmatis] >gpiX87117 MAPGYRA1_1 DNA gyrase [Mycobacterium abscessus] [SUB 75-114]	67	43
60	5	3353	2514	gpJ048361	M.barkeri ATPase alpha and beta subunit (alpha and alphaB) genes, complete cds. [Methanosaerica barkeri] >pirSIA34283 H+-transporting ATP synthase (EC 3.6.1.34) alpha chain - Methanosaerica barkeri	67	51
62	1	572	3	gpIX73141	hemolysin [Serpulina hyodysenteriae]	67	41

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

76	2	2197	677	gpIL2777971	DNA topoisomerase I [Bacillus subtilis]	67	51
91	1	426	gpIM883811	transfer RNA-Leu synthetase [Bacillus subtilis] >pirSIA41882	67	49	
101	9	3437	3616	gpIL207841	leucine-tRNA ligase (EC 6.1.1.4) - Bacillus subtilis	67	37
103	4	3167	4150	gpIU179021	acyl-CoA carboxylase [Cyclotella cryptica] >pirSIA48757 acetyl-CoA carboxylase (EC 6.4.1.2) - Cyclotella cryptica	67	47
113	5	2880	3357	gpIU211921	NrC/NifA-like protein regulator [Escherichia coli]	67	53
121	14	7227	7811	gpIX816421	SecA [Streptomyces lividans]	67	47
135	7	3256	3681	gpIL455211	orf gene product [Wolinella succinogenes] >pirSIS50154 hypothetical protein - Wolinella succinogenes	67	43
					D-alanine permease [Haemophilus influenzae] >gpIU327701HTU32770_3 D alanine permease [Haemophilus influenzae] >gpIU000781HTU00078_28 D alanine permease [Haemophilus influenzae] >gpIU32716HTU32716_5 alanine permease [Haemophilus influenzae] >pirSIS27164099	67	43
136	6	4442	3414	gpIM6311761	helicase [Staphylococcus aureus] >pirSIS27667 DNA helicase pcrA - Staphylococcus aureus >pirSIS39923 pcrA protein - Staphylococcus aureus	67	43
139	1	345	4	gpIZ282011	S.cerevisiae chromosome XI reading frame ORF YKL202w. [Saccharomyces cerevisiae] >pirSIS38038 hypothetical protein YKL201c - yeast (Saccharomyces cerevisiae)	67	41
140	1	1	345	gpIL458931	periplasmic serine protease Do and heat shock protein [Haemophilus influenzae] >gpIU328051HTU32805_12 periplasmic serine protease Do and heat shock protein [Haemophilus influenzae] >gpIU000821HTU00082_28 periplasmic serine protease Do and heat shock protein	67	44
141	1	2	2029	gpIM5586361	Do and heat shock protein pyruvate,orthophosphate dikinase [Zea mays] >gpIS469651S46964S2_1 orthophosphate dikinase [Chloroplast Zea sp.] {SUB 1-154}	67	56
158	4	1690	1217	gpIU425501	A7L gene product [Paramaecium bursaria Chlorella virus 1]	67	39
165	4	922	1485	gpIL086261	abc gene product [Escherichia coli] >gpID353361ECOTSF_24 ATP-binding protein [Escherichia coli]	67	51
213	1	257	48	gpIL384241	polyA polymerase [Bacillus subtilis] >gpIL38424IBACJOJC_7	67	54

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

2	12	9394	9984	gpiD261835I	polyA polymerase [Bacillus subtilis]		
4	11	7319	7885	gpiL4797I	unknown [Bacillus subtilis]	66	39
					methionine aminopeptidase [Bacillus subtilis]	66	42
					>gi D00619 BACSECY_5 B.subtilis secY gene. [Bacillus subtilis]		
					>priSIS0493 methionyl aminopeptidase (EC 3.4.11.18) -		
					Bacillus subtilis		
8	12	6872	6351	gpiM90060I	Streptococcus faecalis H+ ATPase a (atpB),b (atpF),c (atpE),alpha (atpA), beta (atpD) gamma (atpG),delta (atpH),and epsilon (atpC) subunits, complete cds. [Streptococcus faecalis]	66	43
8	34	19445	20911	gpiD14162I	SecY protein [Corynebacterium glutamicum]	66	44
8	44	26260	26472	gpiX67646I	heat-shock protein [Borrelia burgdorferi]	66	43
					>gi M96847 BORGRPEPLS_2 dnaK homologue gene product [Borrelia burgdorferi] >gpiM97912 BORHSP70A_1 70 kDa heat shock protein [Borrelia burgdorferi] >gpiS42385 S42385_1 HSP70 homolog [Borrelia burgdorferi, CA12 isol]		
16	4	2774	1785	gpiZ15056I	mutE gene product [Bacillus subtilis] >priSIB4769I UDP-N-acetylmuramoyl-L-alanyl-D-glutamate-2,6-diaminopimelate ligase (EC 6.3.2.13) - Bacillus subtilis	66	46
16	30	18350	17628	gpiL45659I	hypothetical protein (GB:L10328_69) [Haemophilus influenzae] >gi U32781 HTU32781_4 hypothetical protein (GB:L10328_69) [Haemophilus influenzae] >gi U00079 HTU00079_61 hypothetical protein (GB:L10328_69) [Haemophilus influenzae]	66	48
					>gi U32727 HTU32727_5 imer		
19	1	386	2	gpiL18927I	DNA polymerase III epsilon subunit [Buchnera aphidicola]	66	49
20	1	422	1105	gpiU00039I	E. coli chromosomal region from 76.0 to 81.5 minutes. [Escherichia coli] >priSIS47705 hypothetical protein f648 - Escherichia coli	66	48
21	1	676	1398	gpiJ05534I	Escherichia coli ATP-dependent clp proteolytic component (clpP) gene, complete cds. [Escherichia coli]	66	44
					>priSB36575 ATP-dependent Clp proteinase (EC 3.4.21.-) chain P precursor - Escherichia coli		
21	9	7147	8034	gpiD64006I	hypothetical protein [Synechocystis sp.]	66	48

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

27	4	1615	2550	gpiU249782I	peptide chain release factor 1 [Bacillus subtilis] >pirSIS55437	66	44
28	4	3054	1852	gpiL40822I	peptide chain release factor 1 - Bacillus subtilis	66	50
28	7	3775	3392	gpiL23147I	phosphoglucoisomerase-like protein [Chlamydia trachomatis]	66	40
35	32	18159	17989	gpiU10577I	acetate kinase [Methanoscarcina thermophila] >pirSISB49338 acetate kinase (EC 2.7.2.1) - Methanoscarcina thermophila	66	40
35	34	18290	18144	gpiD43920I	bone sialoprotein II [Gallus gallus]	66	53
					DNA (cytosine-5')-methyltransferase [Gallus gallus]	66	40
					>gpiD43920ICHKMETASE 1 DNA (cytosine-5')-methyltransferase [Gallus gallus]	66	40
39	3	950	777	gpiL47164I	dnrQ gene product [Streptomyces peucetius]	66	46
					>gpiL47164ISTMDNRQ_1 dnrP gene product [Streptomyces peucetius]	66	46
47	12	6842	5379	gpiU18744I	MgeE [Bacillus firmus]	66	39
48	6	3736	3029	gpiL27492I	triosephosphate isomerase [Thermotoga maritima]	66	50
52	6	6983	6273	gpiX73141I	hemolysin [Serpulina hyodysenteriae]	66	40
53	9	5378	3882	gpiM29495I	unknown protein [Synechococcus sp.] >pirSISQ3YCRQ hypothetical protein (recA 3' region) - Synechococcus sp. (PCC 7002) (fragment)	66	35
57	3	1788	2258	gpiJ03218I	S.acidocaldarius membrane-associated ATPase alpha subunit gene, complete cds. [Sulfolobus acidocaldarius] >pirSIA28552 H+-transporting ATP synthase (EC 3.6.1.34) alpha chain, membrane-associated - Sulfolobus acidocaldarius	66	42
57	5	3732	4421	gpiU47274I	AIAOH+ ATPase, subunit D [Methanocarcina mazaeii]	66	39
58	3	1805	1167	gpiU11045I	ORF-C gene product [Buchnera aphidicola]	66	42
61	13	8053	7388	gpiU45426I	heat shock protein HTPG [Borrelia burgdorferi]	66	37
					>gpiL32145IBORHPTG_1 C62.5 heat shock protein [Borrelia burgdorferi] (SUB 497-575)	66	37
62	4	1966	2160	gpiU1983II	SpZ12-1 [Strongyllocentrotus purpuratus]	66	53
63	1	2	979	gpiU27343I	MutL [Bacillus subtilis]	66	42
67	2	2605	1454	pirSIA480I	mucin MG2=low molecular weight salivary glycoprotein - human	66	47
			8		>pirSIS29115 mucin MG2b-T2 - human {SUB 143-168}	66	47
76	1	740	3	gpiU27797I	DNA topoisomerase I [Bacillus subtilis]	66	43

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

121	5	1550	1341	gpiU206691	ribosomal protein S21 [Myxococcus xanthus]	66	45
	8	3037	3657	gpiL097031	penicillin-binding protein [Bacillus subtilis] >gpiZ68230IBSYLLSPO_5 high molecular weight penicillin binding protein [Bacillus subtilis]>pirlSC53292 penicillin-binding protein 2B - Bacillus subtilis>gpiZ225865IBSSPOVD_1 Pbp2B	66	37
127	1	636	4	gpiL448281	Bacillus subtilis 1 [SUB 6 hypothetical protein (SP:P30143) [Haemophilus influenzae] >gpiU32703[HIU32703_7 hypothetical protein (SP:P30143) [Haemophilus influenzae]>pirlU00070[HIU00070_87 hypothetical protein (SP:P30143) [Haemophilus influenzae] >gpiU32812[HIU32812_7 amino acid per	66	57
129	2	1410	2207	gpiU293991	major outer sheath protein [Treponema denticola] isoleucyl-tRNA synthetase [Homo sapiens]	66	36
131	1	471	4	gpiU049531	L20 gene product [Bacillus subtilis]>pirlSIS18439 Ribosomal protein L21 - Bacillus subtilis	66	49
141	3	2103	2453	gpiX595281	mab-21 gene product [Caenorhabditis elegans] Escherichia coli K-12 chromosomal region from 92.8 to 00.1 minutes. [Escherichia coli]>pirlSIS56374 hypothetical protein F342 - Escherichia coli	66	41
158	3	702	463	gpiU198611	RNA-directed RNA polymerase (EC 2.7.7.48) - equine arteritis virus >gpiX534591[OEAV_1 Equine arteritis virus (EAV) RNA genome. [Equine arteritis virus] [SUB 1-1727]	66	66
165	1	190	2	gpiU140031	transcription-repair coupling factor [Bacillus subtilis]	66	52
266	1	86	289	pirlSIRRWV	unknown [Escherichia coli]	66	44
			EV				
2	15	13965	13210	gpiD261851	ribosomal S8 protein [Thermus thermophilus] >pirlSAS53870 ribosomal protein S8 - Thermus aquaticus >pirlSIS51059 ribosomal protein S8 - Thermus aquaticus [SUB 1- 28]	65	41
5	1	74	457	gpiD835361		65	52
8	26	13905	14519	gpiU439291	S3 [Bacillus subtilis]	65	47
8	30	17102	17455	gpiX795511	>pirlSAS53870 ribosomal protein S8 - Thermus aquaticus >pirlSIS51059 ribosomal protein S8 - Thermus aquaticus [SUB 1- 28]	65	59
8	31	17453	18004	pirlSIR5BS0	ribosomal protein L6 - Bacillus stearothermophilus	65	48
16	36	22233	21427	pirlSIA4166	hypothetical protein 1 - Enterococcus faecalis plasmid pCF10	65	40

TABLE 2.
Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				2 (fragment)		
19	3	2367	2894	gpiU18937I histidyl-tRNA synthetase homologue [Homo sapiens] >gpiU18936IHSU18936_1 histidyl-tRNA synthetase homologue [Homo sapiens] {SUB 1-36}	65 42	
19	16	13750	13091	pirSJC1317 protein-methionine-S-oxide reductase (EC 1.8.4.6) - Escherichia coli	65 53	
21	18	11207	10506	gpiU18997I Escherichia coli K-12 chromosomal region from 67.4 to 76.0 minutes. [Escherichia coli]	65 44	
29	26	12129	12341	gpiS79915I His [Drosophila Alcaligenes eutrophus]	65 38	
33	2	372	677	gpiM69036I protein H [Alcaligenes eutrophus] >priSTA38120 phbH protein -	65 42	
35	14	8692	9093	gpiX87899I thioredoxin/thioredoxin reductase hybrid protein [Mycobacterium leprae] >gpiL39923MSGDNAB_14 thioredoxin reductase/thioredoxin [Mycobacterium leprae]	65 42	
35	33	18117	18719	gpiL4607II hypothetical protein (SP:P26242) [Haemophilus influenzae] >gpiU32822IHTU32822_12 hypothetical protein (SP:P26242) [Haemophilus influenzae] >gpiU00084IHTU00084_34 hypothetical protein (SP:P26242) [Haemophilus influenzae]	65 41	
36	18	11816	13033	gpiM20793I S.typhimurium D-alanine:D-alanine ligase (ddlA) gene, complete cds. [Salmonella typhimurium] >pirSICCEBDT D-alanine-D-alanine ligase (EC 6.3.2.4) A - Salmonella typhimurium	65 44	
42	8	3884	4252	gpiX02499I Rhodospirillum rubrum atp operon. [Rhodospirillum rubrum] >pirSIS08579 hypothetical protein 2 - Rhodospirillum rubrum >gpiX02499I RRA_1P_4 Rhodospirillum rubrum atp operon. [Rhodospirillum rubrum] {SUB 592-811}	65 48	
44	6	3647	2583	gpiX95669I thdF gene product [Borrelia burgdorferi] >gpiZ12160IBBGIDAG_1 thdF gene product [Borrelia burgdorferi] {SUB 429-463}	65 50	
47	4	2157	1501	gpiL47709I ypiA gene product [Bacillus subtilis]	65 50	
47	20	12498	12091	gpiL44847I tRNA (guanine-N1)-methyltransferase [Haemophilus influenzae] >gpiU322705IHTU322705_6 tRNA (guanine-N1)-methyltransferase	65 47	

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				[Haemophilus influenzae] >gpiU000711HU00071_17 tRNA (guanine-N1)-methyltransferase [Haemophilus influenzae]		
48	9	6125	7117	gpiU185321	>gpiU32814HTU32814_6 tRNA	
70	3	908	2131	gpiX731241	Bex [Bacillus subtilis] ipa-65d gene product [Bacillus subtilis] >pirSISS39720 hypothetical	
72	1	314	18	gpiU090051	HfIC [Vibrio parahaemolyticus]	
79	5	2430	3350	gpiU437391	FilF [Borrelia burgdorferi] >gpiL76303BORTSA_10 flagellar basal body rod protein [Borrelia burgdorferi]	
81	9	5977	7617	gpiM915931	Mycoplasma mycooides SRPM54 gene, complete cds. [Mycoplasma mycooides] >pirSISS35480 hypothetical protein 1 - Mycoplasma mycooides (SGC3) >pirSISS27590 hypothetical protein 1 - Mycoplasma mycooides (SGC3) (fragment) [SUB 2-422]	
83	13	4268	4939	pirSIIPC230	X-Pro aminopeptidase (EC 3.4.11.9) L13K - guinea pig (fragment) >pirSIIPC2310 X-Pro aminopeptidase (EC 3.4.11.9)	
85	7	3480	2587	gpiM645191	K13K - guinea pig (fragment) [SUB 1-26] transport protein [Escherichia coli] >pirSIIC40840	
88	8	3582	3902	gpiX776361	spermidine/putrescine transmembrane protein C - Escherichia coli putative amino acid binding subunit [Bacillus subtilis]	
				>pirSISS2381 probable amino acid binding protein - Bacillus subtilis		
89	5	2189	2815	gpiM301981	recQ gene product [Escherichia coli]	
98	2	1233	760	gpiX779251	dUTPase [Candida albicans] >pirSISS42871 dUTP pyrophosphatase (EC 3.6.1.23) - yeast (Candida albicans)	
101	1	440	3	gpiU013221	ribonucleotide reductase small subunit [Plasmidum falciparum] >pirSB49412 ribonucleoside-diphosphate reductase (EC 1.17.4.1) small subunit (EC 1.17.4.1) - Plasmidum falciparum	
104	15	5605	6159	gpiM134621	cheW gene product [Escherichia coli] >pirSIQRRECCW chemoaxis protein cheW - Escherichia coli	
115	9	4532	4167	gpiJ054781	alkyl hydroperoxide reductase [Salmonella typhimurium]	
121	1	774	208	gpiX790871	methionyl-tRNA formyltransferase [Thermus aquaticus thermophilus] >pirSI55228 fm1 protein homolog - Thermus	

TABLE 2. *Treponema pallidum* - Putative coding regions of novel proteins similar to known proteins

				aquaticus			
121	18	8556	8110	gpID261851 cell division protein [Bacillus subtilis]			65
125	2	1975	17631	gpIM764421 crystal protein [Bacillus thuringiensis] >pirSI[A41969 crystal protein, 40K - Bacillus thuringiensis >gpIM908431BACCRY_1 crystal protein [Bacillus thuringiensis] [SUB 1-39]			65
144	4	1050	649	gpIX696011 p93 [Borrelia burgdorferi]			46
150	6	4798	5175	gpIL477091 poly(A) polymerase [Bacillus subtilis]			35
150	8	7329	6343	gpID261851 stage 0 sporulation [Bacillus subtilis] >gpIX625391BSORIGS_11 spoQ93 gene product [Bacillus subtilis] >pirSIA38536 spo0193 protein - Bacillus subtilis			50
151	1	1231	2	gpIM302971 B.subtilis recombination and sporulation protein (recN, spoIVB) genes, complete cds, arginine hydroximase resistance (ahcC) gene, 3' end [Bacillus subtilis] >pirSB35128 recN homolog - Bacillus subtilis			47
1	7	2258	3169	gpID833361 proline-tRNA ligase [Escherichia coli]			37
1	10	3063	3620	gpIL251051 aminocycl-tRNA synthetase [Chlamydia trachomatis]			65
2	5	2357	4033	gpIL152021 Bacillus subtilis comE operon encoding ORF1, ORF2, ORF3 and Reverse-ORF genes, complete cds. [Bacillus subtilis] >pirSS39865 ComE ORF3 - Bacillus subtilis			41
4	7	4689	5108	gpIU394831 EMG2 [Escherichia coli]			31
5	3	1285	2172	gpID640001 hypothetical protein [Synechocystis sp.]			46
8	40	23622	24191	gpID2853501 heat shock protein GrpE homolog [Synechococcus sp.] >pirSPC2235 GrpE protein - Synechococcus sp. (PCC 7942), (fragment)			43
8	42	23945	24448	gpIM849641 heat shock protein [Bacillus subtilis] >gpIX514771BSGRPE_1 Bacillus subtilis DNA for grpE gene and dnaK gene (partial). [Bacillus subtilis] >pirSIS08418 heat shock protein grpE - Bacillus subtilis			38
11	19	8921	9778	gpIU293991 major outer sheath protein [Treponema dentitcola]			38
16	20	13266	12769	gpIM85240 flagellar protein [Escherichia coli]			33
16	24	15274	14606	gpIZ500981 pentose-5-phosphate-3-epimerase [Solanum tuberosum]			44
37	10	4141	4788	gpIU189971 Escherichia coli K-12 chromosomal region from 67.4 to 76.0			40

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

41	14	5354	5845	gplD641161	Stb [Bacillus subtilis] >gplD49781IBACSRBA_2 Stb [Bacillus subtilis] >pirSIC4093 signal recognition particle receptor alpha chain homolog - Bacillus subtilis	64	46
47	11	5268	5486	pirSIA4070	tenascin-X precursor - human >gplX71937HSXBKVIII_1 fibrinogen [Homo sapiens] {SUB 3356-3566} >pirSIC42175 tenascin homolog 3.9kF3.1 - human (fragment) {SUB 1849-1936}	64	42
54	27	8658	8260	gplM609171	purine nucleoside phosphorylase [Escherichia coli] >gplU140031ECO[W93_295 purine-nucleoside phosphorylase [Escherichia coli]] >pirSIA27854 purine-nucleoside phosphorylase (EC 2.4.2.1) - Escherichia coli	64	51
58	1	575	3	gplX840191	orf3 gene product [Zymomonas mobilis] >gplX840191ZMDNAGRP_3 orf3 gene product [Zymomonas mobilis]	64	36
62	10	4019	4354	gplL459801	spermidine/putrescine transport ATP-binding protein [Haemophilus influenzae] >gplU32813IHTU32813_12 spermidine/putrescine transport ATP-binding protein [Haemophilus influenzae] >gplU00083IHTU00083_23 spermidine/putrescine transport ATP-binding protein [Haemophilus influenzae]	64	49
64	21	9281	9610	gplZ542381	T28C6.1 [Caenorhabditis elegans]	64	52
70	23	15905	15985	gplD261851	cysteinyl-tRNA synthetase [Bacillus subtilis] >gplJ14580IBACGLUSYN_6 cysteinyl-tRNA synthetase [Bacillus subtilis] >gplX73989IBSCTS_1 cysteine-tRNA ligase [Bacillus subtilis] >pirSIC53402 cysteine-tRNA ligase (EC 6.1.1.16) - Bacillus subtilis	64	50
88	4	2017	2583	gplU185391	FliY [Escherichia coli]	64	33
100	42	18967	20169	gplX816421	orf gene product [Wolinella succinogenes] >pirSIS50154 hypothetical protein - Wolinella succinogenes	64	48
112	15	10532	9804	gplU308211	putative protein of 244 amino acids [Cyanelle Cyanophora paradoxa]	64	42
115	7	2550	3152	gplX916551	lepA gene product [Bacillus subtilis] >gplD17650IBACGPR_4	64	45

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				ORF80 protein [Bacillus subtilis] [SUB 1-327]		
121	27	15390	17069	gplI448761 ATP-dependent RNA helicase [Haemophilus influenzae] >gplU327091 HIU32709_6 ATP-dependent RNA helicase [Haemophilus influenzae]>gplU00071 HIU00071_47 ATP- dependent RNA helicase [Haemophilus influenzae] >gplU328181 HIU32818_6 RNA (?) helicase [Haemophilus in	64	47
				major outer sheath protein [Treponema denticola]	64	51
123	1	44	439	gplU293991 major outer sheath protein [Treponema denticola]	64	38
129	1	1208	3	gplU293991 major outer sheath protein [Treponema denticola]	64	42
170	1	874	92	gpmM548841 xpb gene product [Escherichia coli]>gplU2875 ECU2875_44 site-specific integrase/recombinase, with xerC [Escherichia coli] >pirlSA39202 recombinase XerD - Escherichia coli	64	42
552	1	3	407	gpmM278691 B.subtilis ahnC gene, encoding an arginine repressor/activator protein. [Bacillus subtilis]	64	39
636	2	265	642	gpmM264141 RNA polymerase alpha-core-subunit [Bacillus subtilis] >gplL4797 BACRPLP_21 RNA polymerase alpha-core-subunit [Bacillus subtilis]>pirlSE32307 DNA-directed RNA polymerase (EC 2.7.7.6) alpha chain - Bacillus subtilis >gplM13957 BACRPOA_3 B.subtilis DNA se	64	40
				limb deformity protein [Gallus gallus]>pmlS38780 limb deformity protein - chicken	64	64
685	1	2	286	gpxX626811	64	64
693	2	383	168	gplU410471 Genesis [Mus musculus]>gplL13202 RATHFH2_1 HNF-3/fork- head homolog-2 [Rattus norvegicus] [SUB 129-229]	64	58
2	9	5368	6384	gplI448211 hypothetical protein (SP:P33643) [Haemophilus influenzae] >gplU32702 HIU32702_15 hypothetical protein (SP:P33643) [Haemophilus influenzae]>gplU00070 HIU00070_80 hypothetical protein (SP:P33643) [Haemophilus influenzae] >gplU32811 HIU32811_16 pseudoU synt	63	38
8	25	13565	13933	gplVZ216771 ribosomal protein L22 [Thermotoga maritima]>priSIS40193	63	45
				ribosomal protein L22 - Thermotoga maritima	63	45
8	47	27536	27841	gplD640061 hypothetical protein [Synechocystis sp.]	63	36
16	29	17481	16768	gplD261851 unknown [Bacillus subtilis]>gplX62539 BSORIGS_4_B.subtilis genes rpmH, rpmA, 50kd, gida and gidB. [Bacillus subtilis]	63	39

TABLE 2. *Treponema pallidum* - Putative coding regions of novel proteins similar to known proteins

				>gp Z14225 [SRNPASPO_3 Jag [Bacillus subtilis]]>pintSIS18074
21	12	8978	9610	gp U034971 valy tRNA synthetase [Escherichia coli]
22	17	11185	11997	gp U295801 mazG gene product [Escherichia coli]
22	19	13409	13194	gp L463191 endonuclease III [Haemophilus influenzae] >gp U32842 HTU32842_1 endonuclease III [Haemophilus influenzae] >gp U00861 HTU0086_51 endonuclease III [Haemophilus influenzae] >gp U32781 HTU32788_7 endonuclease III [Haemophilus influenzae] >gp U32781 HTU32788_7 endonuclease III [Haemophilus influenzae] >pir S1G64136 endonu
29	22	10192	9176	gp U293991 major outer sheath protein [Treponema denticola]
39	9	4947	3727	gp U26401 galactokinase [Homo sapiens] >gp U26401 HTSU26401_1 galactokinase [Homo sapiens]
42	9	4123	4872	gp L449831 primosomal protein replication factor [Haemophilus influenzae] >gp U32718 HTU32718_10 primosomal protein replication factor [Haemophilus influenzae] >gp U00072 HTU00072_68 primosomal protein replication factor [Haemophilus influenzae] >gp U32827 HTU32827_
47	26	15583	14249	gp Z150561 murD gene product [Bacillus subtilis] >pintSID47691 UDP-N-acetyl muramoylalanine-D-glutamate ligase (EC 6.3.2.9) - Bacillus subtilis
49	8	6603	5830	gp L450451 hypothetical protein (GB:D26185_130) [Haemophilus influenzae] >gp U32724 HTU32724_5 hypothetical protein (GB:D26185_130) [Haemophilus influenzae] >gp U00073 HTU00073_26 hypothetical protein (GB:D26185_130) [Haemophilus influenzae] >gp U32832 HTU32832_13_A
54	18	6674	6345	pintSIB6122 collagen alpha 1(IV) chain - rabbit (fragment)
56	8	3790	4077	gp D451631 embryonic muscle myosin heavy chain [Halocynthia roretzi]
58	4	2308	1829	gp L461861 H. influenzae predicted coding region HTI1555 [Haemophilus influenzae] >gp U32830 HTU32830_1 H. influenzae predicted coding region HTI1555 [Haemophilus influenzae] >gp U0085 HTU00085_26 H. influenzae predicted coding region

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

					HU15353 [Haemophilus influenzae]		
61	12	6015	6527	gpiU189971	gipG gene product [Escherichia coli]	63	32
69	13	7256	7561	gpiX693771	glycine-rich RNA-binding protein [Arabidopsis thaliana] >pisS31443 glycine-rich RNA-binding protein (clone A81) - Arabidopsis thaliana (fragment) >gpiZ181891ATTSG0693_1	63	41
					GLYCINE-RICH RNA-BINDING PROTEIN [Arabidopsis thaliana] [SUB 1-66]		
77	2	904	1065	pirSTNLJS ₂	trans-activating transcriptional regulatory protein - simian immunodeficiency virus SV40m (type 3, isolate STL V-3agm)	63	47
81	19	16307	17188	gpIM870491	uridine phosphorylase [Escherichia coli] >gpiX156891ECUDP_1 E. coli udp gene for uridine phosphorylase (EC 2.4.2.3). [Escherichia coli] >pisSIS05491 uridine phosphorylase (EC 2.4.2.3) - Escherichia coli	63	48
88	1	1	966	gpiLA54391	ribosomal protein S4 [Haemophilus influenzae] >gpiU32762HU32762_15 ribosomal protein S4 [Haemophilus influenzae] >gpiU00077HU00077_61 ribosomal protein S4 [Haemophilus influenzae] >gpiU327081HU32708_16 ribosomal protein S4 [Haemophilus influenzae] >P	63	40
111	2	1345	1902	gpiUI2513	PPI-dependent phosphofructo-1-kinase [Entamoeba histolytica] >gpiU12513EHU12513_1 PPI-dependent phosphofructo-1-kinase [Entamoeba histolytica] >pisSIS52082 PPi-dependent phosphofructo-1-kinase - Entamoeba histolytica	63	52
121	4	1277	1113	gpiI452631	poly peptide deformylase (formylmethionine deformylase) [Haemophilus influenzae] >gpiU32745HU32745_1 polypeptide deformylase (formylmethionine deformylase) [Haemophilus influenzae] >gpiU000751HU00075_42 polypeptide deformylase (formylmethionine deformylase)	63	44
167	3	1293	1490	gpiD265621	'ribosome releasing factor' [Escherichia coli] >gpiJ05113ECORRFX_1 E.coli ribosome releasing factor gene, complete cds. [Escherichia coli] >gpiD13334IECOSMBA_3 ribosome releasing factor (FRR) [Escherichia coli] >gpd83536IECOTSF_3 ribosome releasing fact	63	42
646	2	417	674	gpiD378501	core, env, and part of E2/NS1 [Hepatitis C virus]	63	47

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				>gpi D37850 HPCORE13_1 core, env and part of E2/NS1
5	5	2360	3250	gpi D83536 [Hepatitis C virus]
10	18	9417	9190	gpi U26536 unknown [Escherichia coli]
29	6	2643	2431	glycerophosphoryl diester phosphodiesterase [Escherichia coli]
38	3	631	1548	YOR3174c gene product [Saccharomyces cerevisiae]
42	1	999	679	T25B9.9 [Caenorhabditis elegans]
42	13	7604	7918	GlnQ [Mycoplasma pneumoniae]
44	24	14444	13665	transcription factor [Thermus thermophilus] >pir SIA4528
45	13	10253	11362	transketolase [Thermus thermophilus] >SUB_61_244
48	19	12383	11879	transketolase [Xanthobacter flavus] >gpi U33064 XF-U33064_1
55	1	2	475	transketolase [Xanthobacter flavus]
62	3	1895	945	hypothetical protein (SP:P33943) >gpi U32702 HTU32702_11 hypothetical protein (SP:P33943)
64	4	1781	2050	[Haemophilus influenzae] >gpi U00070 HTU00070_76 hypothetical protein (SP:P33943) [Haemophilus influenzae]
69	5	2474	1995	>gpi U10397 galactose binding protein [Citrobacter freundii] >pir SIS15354 galactose-binding protein - Citrobacter freundii
69	12	6034	6690	YHR155w gene product [Saccharomyces cerevisiae] >pir SIS46754 hypothetical protein YHR155w - yeast
70	21	14901	15395	(Saccharomyces cerevisiae)
				>pir SIS46754 hypothetical protein [Synechocystis sp.]
				merozoite surface antigen 1 precursor - Plasmodium vivax
				histidine rich protein [Escherichia coli] >gpi L28082 ECOSLYD_1
				styD gene product [Escherichia coli] >gpi L3261 ECOSLYDX_2
				styD gene product [Escherichia coli] >gpi U18997 ECOUW67_273
				histidine rich protein [Escherichia coli] >pir SIA49987 metal-binding pro
				fragment)

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

83	8	3198	3713	gpID2543281	unknown [Schizosaccharomyces pombe]		62	44
87	1	1678	2	gpIS57681	EF-G [Thermotoga maritima]		62	44
89	13	11242	11454	gpIU170131	Oct1 [Rattus norvegicus]		62	31
91	3	902	1405	gpIM885811	transfer RNA-Leu synthetase [Bacillus subtilis] >pirISIA41882		62	46
94	5	2441	2160	gpID104831	ORF [Escherichia coli] >pirISIQQECAF1 hypothetical 38.8K protein (ftsI 5' region) - Escherichia coli >gpiX55034 EC2MIN_9 E. coli 2 minute region. [Escherichia coli] [SUB 34-346]		62	43
95	3	1626	277	gpIU293991	major outer sheath protein [Treponema denticola]		62	43
102	5	3413	2340	gpIX021641	E.coli ponA gene for penicillin-binding protein 1A (PBP 1A). [Escherichia coli] >pirISIZPECPA penicillin-binding protein 1A - Escherichia coli		62	40
104	19	7930	8310	gpIU368401	Escherichia coli K-12 genome, approximately 57 minutes. [Escherichia coli]		62	46
110	1	16	429	gpIX644511	gcpE gene product [Escherichia coli] >pirIS23058 gcpE protein - Escherichia coli		62	51
111	1	698	1003	gpID2266531	laminin M chain (merosin) [Homo sapiens]		62	33
115	1	803	231	gpID261851	unknown [Bacillus subtilis] >gpIL14580 BACGLUSYN_3 unknown [Bacillus subtilis] >pirISIA53402 glutamate-t-RNA ligase (EC 6.1.1.17) - Bacillus subtilis (fragment) [SUB 137-158]		62	47
					Synthetic Clostridium MP flavodoxin gene, complete cds. [Artificial gene] >pirISFXCLEX flavodoxin - Clostridium sp.		62	47
116	7	2360	1719	gpU050151	Neisseria meningitidis dTDP-D-glucose 4,6-dehydratase (rfbB), glucose-1-phosphate thymidylyl transferase (rfbA) and rfbC genes, complete cds and UDP-glucose-4-epimerase (gale) pseudogene. [Neisseria meningitidis]		62	39
141	5	3195	3944	gpIM245371	GTP-binding protein [Bacillus subtilis] >pirISIB32804 GTP-binding protein, spoOB 3'-region - Bacillus subtilis >gplK026661BACSP00B2_2 Bacillus subtilis spoOB early sporulation gene, complete cds. [Bacillus subtilis] [SUB 1-65]		62	48
152	1	3	1964	gpJ032941	uvrB gene product [Bacillus subtilis] >gplJ032941BACAPKU_2 deoxyribodipyrimidine photolyase [Bacillus subtilis]		62	48

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				>pirSIS437192_uvrB protein - <i>Bacillus subtilis</i>		
160	2	272	874	gpi Z121601 gldA gene product [Borrelia burgdorferi] >gpi Z121601BBGIDAG_1 division protein [Borrelia burgdorferi] (SUB_529_593) >gpi X95669BBTIDFGID_2 gldA gene product [Borrelia burgdorferi] (SUB_1-29) >gpi X95668BBGIDMOXR_1 gldA gene product [Borrelia burgdorferi]	62	34
194	1	378	4	gpi U049531 isoleucyl-tRNA synthetase [Homo sapiens]	62	44
574	1	439	20	gpi D261851 unknown [Bacillus subtilis] >gpi U02604IBSU02604_3 ORF Y [Bacillus subtilis] (SUB_1-260)	62	40
599	1	1	321	gpi D261851 unknown [Bacillus subtilis]	62	37
645	2	592	200	gpi U293991 major outer sheath protein [Treponema dentitcola] methionine aminopeptidase [Bacillus subtilis] >gpi D00619BACSECY_5 B.subtilis secY gene. [Bacillus subtilis] >pirSIS0493 methionyl aminopeptidase (EC 3.4.11.18) - <i>Bacillus subtilis</i>	62	37
4	12	7734	8114	gpi Z479711 ribosomal protein L23 [Thermotoga maritima] >pirSIS40190 ribosomal protein L23 - Thermotoga maritima	61	47
8	22	11941	12378	gpi Z216771 C09G5.8 [Caenorhabditis elegans]	61	41
13	12	7268	7047	gpi Z467921 rodA gene product [Escherichia coli] >pirSISBVECRD rod shape-	61	38
30	17	12484	13842	gpi M228571 determining protein mrdB - Escherichia coli	61	39
35	30	17342	18001	pirSIS156181E4 protein - human papillomavirus type 2a	61	47
35	41	24978	24346	gpi U510321 D9651.10 gene product [Saccharomyces cerevisiae]	61	38
44	14	8614	7922	gpi A21151 insulin-activated amino acid transporter [<i>Mus musculus</i>]	61	38
49	12	86688	8234	gpi D261851 >pirSISIC4149 adipocyte amino acid transporter - mouse	61	46
57	8	7197	6619	gpi D640011 hypothetical protein [Synechocystis sp.]	61	37
63	27	16398	14845	gpi X731241 ipa-68d gene product [Bacillus subtilis] >pirSIS39723 hypothetical protein - <i>Bacillus subtilis</i>	61	44
65	2	396	214	gpi L147451 C.elegans cosmid C02F5. [Caenorhabditis elegans] >pirSIS44608 C02F5.6 protein - <i>Caenorhabditis elegans</i>	61	44
79	1	2	895	gpi X784781 sss gene product [Pseudomonas aeruginosa] >pirSIS43156 sss protein - <i>Pseudomonas aeruginosa</i>	61	49

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

94	2	1580	750	epIX624371	[UDP-N-acetyl]muramoylalanyl-D-glutamyl-2, 6-diaminopimelate--D-alanyl-D-alanine ligase [Synechocystis sp.] >pirSIS49610		61	43
					UDP-N-acetyl]muramoylalanyl-D-glutamyl-2, 6-diaminopimelate--D-alanyl-D-alanine ligase (EC 6.3.2.15) - Synechocystis sp. (PCC 6803)			
102	1	934	101	gpIL458491	single-stranded-DNA-specific exonuclease [Haemophilus influenzae] >gpiU32801IHU32801_1 single-stranded-DNA-specific exonuclease [Haemophilus influenzae]		61	43
					>gpiU00081IHU00081_78 single-stranded-DNA-specific exonuclease [Haemophilus influenzae] >gpiU32746IH			
103	2	693	1931	gpID261851	similar to B. subtilis DnaH [Bacillus subtilis]		61	38
104	12	4994	54044	gpiU2833771	Escherichia coli K-12 genome; approximately 65 to 68 minutes. [Escherichia coli]		61	48
116	13	5011	4265	gpiU328471	ribosomal protein S1 homolog, RNA-binding protein. [Haemophilus influenzae]		61	39
117	5	3484	2504	gpM340661	trigger factor [Escherichia coli] >pirSIS1A36129 trigger factor protein - Escherichia coli >gpiU05334IECOCCLPPA_1 Escherichia coli ATP-dependent clp protease proteolytic component (clpP) gene, complete cds. [Escherichia coli] [SUB 390-432]		61	29
121	16	8072	7830	gpIZ470471	unknown [Saccharomyces cerevisiae] >gpiUZ46902ISCS8224_5		61	19
128	1	1415	135	gpIL260511	UDP-N-acetylglucosamine 1-carboxyvinyl transferase [Acinetobacter calcoaceticus]		61	39
142	5	7013	8107	gpID261851	DNA polymerase III subunit [Bacillus subtilis]		61	43
					>gpiX17014BSRECM_1 dnaZX gene product [Bacillus subtilis]			
					>pirSIS13786 DNA-directed DNA polymerase (EC 2.7.7.7) III chain dnaX - Bacillus subtilis >gpiX06803BSDNAZZX_1 Bacillus subtilis DNA for dnaZX-like O			
150	7	5018	6346	gpIL477091	poly(A) polymerase [Bacillus subtilis]		61	46
324	1	1	192	gpIL455211	D-alanine permease [Haemophilus influenzae]		61	42
					>gpiU32770IHU32770_3 D-alanine permease [Haemophilus influenzae] >gpiU00078IHU00078_28 D-alanine permease [Haemophilus influenzae] >gpiU32716IHU32716_5 alanine			

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				permease [Haemophilus influenzae] >pirSISIH64099		
328	2	308	296	gplX665041 Exel gene product [Aeromonas hydrophila] >pirSIE49905 protein	61	38
681	2	239	466	secretion operon exo protein L - Aeromonas hydrophila	61	33
8	37	22115	22876	gplM264141 Thermus thermophilus fus gene for elongation factor G (EF-G). [Thermus aquaticus thermophilus] >pirSIEFIWG translation elongation factor G - Thermus aquaticus	60	38
				>gpl479711BACRPLP_21 RNA polymerase alpha-core-subunit [Bacillus subtilis]		
				>pirSIE32207 DNA-directed RNA polymerase (EC 2.7.7.6) alpha chain - Bacillus subtilis		
				>gplM139571BACRPOA_3 B subtilis DNA se		
16	41	24204	23581	gplX820711 orf3 gene product [Pseudomonas aeruginosa] >pirSIS40376 hypothetical protein 3 - Pseudomonas aeruginosa	60	26
16	48	28011	28598	gplSIS19739 integral membrane protein - Rhodobacter capsulatus	60	35
21	10	7764	7165	gplL447441 H. influenzae predicted coding region HJ0100 [Haemophilus influenzae] >gplU32695HTU32695 7 H. influenzae predicted coding region HJ0100 [Haemophilus influenzae] >gplU000701HTU00070_7 H. influenzae predicted coding region HJ0100 [Haemophilus influenzae] >	60	34
26	24	13879	14238	gplM59441 mgIA gene product [Escherichia coli] >pirSIB37277_50K membrane-associated protein mgIA - Escherichia coli	60	44
35	7	4135	4671	gplD641161 ORF3 [Bacillus subtilis]	60	41
36	13	8256	7405	gplD504171 ATEC3 [Mus musculus]	60	39
41	3	915	1574	gplM627851 multiphosphoryl transfer protein [Rhodobacter capsulatus] >gplX531501RCFRUOP_1 multiphosphoryl transfer protein [Rhodobacter capsulatus] >pirSIS10639 fruB protein - Rhodobacter capsulatus	60	45
44	7	4075	3578	gplX956691 thdF gene product [Borrelia burgdorferi] >gplZ121601BBGIDAG_1 thdF gene product [Borrelia burgdorferi] [SUB 429-463]	60	41
45	5	6624	5203	E.coli dacC gene for penicillin-binding protein 6. [Escherichia coli] >pirSIB28536 penicillin-binding protein 6 precursor - Escherichia	60	49

TABLE 2.
Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				coli			
46	2	1235	858	gpiU452851 specific 116-kDa vacuolar proton pump subunit [Homo sapiens]	60	40	
47	22	12967	12641	gpiU3973II tRNA (guanine-N1)-methyltransferase [Mycoplasma genitalium]	60	42	
				>pirSIIB64249 tRNA (guanine-N1-)methyltransferase (EC 2.1.1.31) - Mycoplasma genitalium (SGC3)			
70	6	4740	5597	gpiM33977I Dps/ery gene product [Drosophila pseudoobscura] >pirSIA31946 xanthine dehydrogenase (EC 1.1.1.204) - fruit fly (Drosophila pseudoobscura)	60	42	
101	8	3562	3326	gpiM74329I pcx gene product [Drosophila melanogaster] >gpiM25662 DROPEC_I pcx gene product [Drosophila melanogaster] [SUB 546-2483]	60	46	
101	16	8112	7486	gpiD64006I hypothetical protein [Synechocystis sp.]	60	32	
108	1	2	583	gpiZ23080I major vegetative sigma factor [Clostridium acetobutylicum] >pirSIS34307 DNA-directed RNA polymerase (EC 2.7.7.6)	60	38	
				sigma factor sigA - Clostridium acetobutylicum			
124	1	917	3	gpiX75568I ICFG [unidentified] >pirSIS48034 icfG protein - Synechocystis sp. (PCC 6803)>pirSIS38573 ICFG protein - Synechocystis sp. (strain PCC6803)	60	33	
130	1	111	719	gpiL09228I Bacillus subtilis spoVA to sea region. [Bacillus subtilis]	60	47	
				>pirSIS45549 hypothetical protein X7 - Bacillus subtilis			
167	1	1	243	gpiU35149I putative glutamate and asparagine rich protein [Plasmodium chabaudi]	60	43	
175	1	253	65	pirSIA2966 keratin, 65K type II cytoskeletal - human >gpiX05418HSKER65A_1 keratin type II (AA1-215) [Homo sapiens] [SUB 1-215]	60	46	
177	1	1077	199	gpiU43739I FtsW [Borrelia burgdorferi] >gpiX96432BBMRAYFTS_2 ftsW	60	39	
				gene product [Borrelia burgdorferi] [SUB 1-27]			
185	1	524	21	gpiU17010I NADH dehydrogenase, subunit 5 [Mitochondrion Allomyces macrogynus] >gpiU41288IAMU41288_8 NADH dehydrogenase, subunit 5 [Mitochondrion Allomyces macrogynus]	60	45	
570	1	2	532	gpiJ03901I Maize pyruvate,orthophosphate dikinase mRNA, complete cds. [Zea mays]	60	31	

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

582	1	2	376	gpiLU000081	yok [Escherichia coli]		60	40
586	1	265	2	gplL454971	ATP-dependent protease binding subunit [Haemophilus influenzae] >gpiU32767 HTU32767_7 ATP-dependent protease binding subunit [Haemophilus influenzae]		60	37
					>gpiU0078 HTU00078_6 ATP-dependent protease binding subunit [Haemophilus influenzae] >gpiU32713 HTU32713			
606	1	368	3	gpiU231631	OrfUU [Escherichia coli]		60	37
10	20	11015	10356	gpiD500641	PgsA [Bacillus subtilis]		59	38
12	8	4396	4947	gpiU295801	Escherichia coli K-12 genome; approximately 62 minute region. [Escherichia coli]		59	38
17	4	2678	2068	gpiLU000211	u0247g [Mycobacterium leprae]		59	42
19	13	10882	10502	gpiM884891	nonamer binding protein [Mus musculus] >priSIUC4236 V(D)J recombinational signal sequence-dependent DNA joining protein 2 - mouse		59	38
19	19	14336	14151	gpiL493361	tryptophanyl-tRNA synthetase [Clostridium longisporum]		59	40
19	20	15203	14334	gpiL493361	tryptophanyl-tRNA synthetase [Clostridium longisporum]		59	39
30	2	1175	14022	gpiU187921	glutamine-dependent carboxyl phosphate synthase [Babesia bovis]		59	44
30	7	3181	4278	gpiZ223080	primase [Clostridium acetobutylicum] >priSI334306 DNA primase - Clostridium acetobutylicum		59	38
30	8	4242	4493	gpiU098251	acid finger protein [Homo sapiens] >gpiU098251HSU09825_1 acid finger protein [Homo sapiens]		59	37
34	4	2138	2740	gpiL444631	hypothetical protein (GB:U00019_14) [Haemophilus influenzae] >gpiU32687 HTU32687_8 hypothetical protein (GB:U00019_14) [Haemophilus influenzae] >priSI00069 HTU00069_17 hypothetical protein (GB:U00019_14) [Haemophilus influenzae] >gpiU32796 HTU32796_11 oxid		59	37
35	29	17973	17338	gpiZ344691	class II metallothionein with homology to wheat Ec [Zea mays] >gpiU10696 ZMU10696_1 Ec metallothionein class II protein [Zea mays] >priSI47158 metallothionein II - maize		59	54
42	7	3528	4172	gpiM332931	E.coli primosomal protein n' (priA) gene, complete cds, and cytR gene, 5' end. [Escherichia coli]		59	43

TABLE 2.
Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

48	12	9614	8892	gplI449581	Holliday junction DNA helicase [Haemophilus influenzae] >gplU32716[H TU32716_14 Holliday junction DNA helicase [Haemophilus influenzae]>gplU00072[H TU00072_43 Holliday junction DNA helicase [Haemophilus influenzae] >gplU32825[H TU32825_13 Holliday junction	59	40
49	9	7040	6501	gplD261851	unknown [Bacillus subtilis]	59	37
53	3	1302	424	gplI449841	hypothetical protein [SP:P32049] [Haemophilus influenzae] >gplU32719[H TU32719_1 hypothetical protein [SP:P32049] [Haemophilus influenzae]>gplU00072[H TU00072_69 hypothetical protein [SP:P32049] [Haemophilus influenzae] >gplU32828[H TU32828_1 SAM-dependent	59	41
63	25	14385	13177	gplY005441	Escherichia coli gyrA gene, orfX and orfY. [Escherichia coli] >gplM87509[ECOUBIG_1 ubiquinone synthesis-related protein [Escherichia coli]>pirSIA47682 2-octaprenyl-3-methyl-5- hydroxy-6-methoxy-1,4-benzoquinone methyltransferase, UbiG - Escherichia coli	59	40
63	32	19064	18459	gplL195211	Synechococcus sp. four ORFs, three complete, and one 3' end. [Synechococcus sp.]	59	41
64	5	3262	2102	gplDS06171	Cytoplasmic phenylalanyl-tRNA synthetase beta chain [Mitochondrion Saccharomyces cerevisiae] >gplD50617[YSCCHRVIN_47 Cytoplasmic phenylalanyl-tRNA synthetase beta chain [Saccharomyces cerevisiae]>pirS1YFBYAC phenylalanine-tRNA ligase (EC 6.1.1.20) beta	59	40
70	26	16844	17113	gplU245691	sigma factor [Pseudomonas aeruginosa]>gplJ36379[PSEALGT_2 alternative sigma factor [Pseudomonas aeruginosa] >gplL14760[PSEALGUA_1 algU gene product [Pseudomonas aeruginosa]>gplJ14761[PSEALGUB_1 algU gene product [Pseudomonas aeruginosa]>gplU49151[PAU49 smf - Escherichia coli	59	40
76	3	3089	2139	pirSIB4969		59	37
79	8	3348	3995	gplI405011	flagellar MS-ring protein [Borrelia burgdorferi]	59	33
79	10	5190	6212	gplI405021	flagellar export protein [Borrelia burgdorferi] >gplL76303[BORFTSA_12 flagellar export apparatus [Borrelia	59	29

TABLE 2.
Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				[burgdorferi]			
85	8	3841	3398	epIM9239I [transmembrane protein [Escherichia coli] >pirsIC45313 putrescine transport protein porH - Escherichia coli]	59	36	
92	1	1	612	gpIL25603I [protease [Treponema denticola]]	59	33	
94	1	825	118	gpIX55034I [UDP-MurNac-pentapeptide presynthetase (AA -20 to 432) [Escherichia coli] >gplX15432I[ECMURF_1 UDP-MurNAC-pentapeptide presynthetase (AA -20 to 432) [Escherichia coli] >gplD10483I[ECO110K_66 UDP-N-acetylMuramoylalanyl-D-glutamyl-2, 6-diaminopimelate-D-alanyl]	59	43	
106	2	656	1174	gpIJ02774I [Chicken embryonic myosin heavy chain gene, complete cds. [Gallus gallus] >pirsIB24124 myosin heavy chain, EF1W1 - chicken (fragment) [SUB 1-168]	59	20	
114	1	3	272	gpID64000I [hypothetical protein [Synechocystis sp.] >gplU38915ISSU38915_1 LytB [Synechocystis sp. J [SUB 28-406]	59	35	
114	4	995	1462	gpIJ15180I [trA [Mycobacterium leprae]]	59	45	
122	3	1795	2625	gpIX56678I [dcIAE gene product [Bacillus subtilis]]	59	41	
136	5	2556	3410	gpIL46086I [H. influenzae predicted coding region H11454 [Haemophilus influenzae] >gplU32823IHU32823_13 H. influenzae predicted coding region H11454 [Haemophilus influenzae] >gplU00084IHUJ00084_49 H. influenzae predicted coding region H11454 [Haemophilus influenzae]]	59	37	
138	1	609	382	gpID26185I [unknown [Bacillus subtilis]]	59	35	
142	1	231	1328	gpU32164I [NAD(P)H-dependent dihydroxyacetone-phosphate reductase [Bacillus subtilis]]	59	41	
202	1	722	183	gpIL44828I [hypothetical protein (SP:P30143) [Haemophilus influenzae] >gplU32703IHU32703_7 hypothetical protein (SP:P30143) [Haemophilus influenzae] >gplU00070IHU00070_87 hypothetical protein (SP:P30143) [Haemophilus influenzae] >gplU32812IHU32812_7 amino acid per flagellar protein [Escherichia coli]]	59	40	
595	1	1	378	gpIM85240I [nestin [Homo sapiens] >pirsIS21424 nestin - human]	59	30	
4	5	3965	4324	gpIX65964I [nestin [Homo sapiens] >pirsIS21424 nestin - human]	58	36	

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

5	9	6219	6545	gpIZ351[17]	alpha 1,4-glucan phosphorylase type H [Vicia faba]		58	52
9	5	4201	5055	gpIX5493[31]	SPase I [Salmonella typhimurium] >pirSIS12020 signal peptidase I - Salmonella typhimurium		58	41
10	1	2596	1991	gpIX8217[4]	deoxyribose-phosphate aldolase [Bacillus subtilis] >pirSIS4945[5]		58	44
18	4	5354	3198	gpIU4373[91]	deoxyribose-phosphate aldolase (EC 4.1.2.4) - Bacillus subtilis		58	37
					Borreia burgdorferi fesmid clone 31, complete sequence. [Borrelia burgdorferi]			
25	2	1075	149	gpIL1443[7]	flagellar hook-filament junction protein [Bacillus subtilis]		58	34
26	5	3419	3730	gpID2618[51]	unknown [Bacillus subtilis]		58	36
28	19	8985	8803	gpIV0014[1]	Cauliflower mosaic virus genome. [Cauliflower mosaic virus] >pirSISQCV6[5] hypothetical protein 6 - cauliflower mosaic virus (strain Strasbourg)		58	58
29	20	9207	8524	gpU2939[91]	major outer sheath protein [Treponema denticola]		58	45
41	2	547	945	gpIM219[94]	E.coli cysK gene, 3' end, ptsH, ptsJ, and crt phototransferase system genes, complete cds. [Escherichia coli]		58	37
					>pirI02796[ECOPTSII_2] ptsI gene product [Escherichia coli]			
					>pirSIWQECPI phosphotransferase system enzyme I (EC 2.7.3.9) - Escherichia coli			
45	10	8472	9128	pirSIB550[5]	endothelial monocyte-activating protein II precursor - human		58	27
48	21	13329	12793	gpIU1899[7]	Escherichia coli K-12 chromosomal region from 67.4 to 76.0 minutes. [Escherichia coli] >gpI1566[IECU1566]_1 HhoA [Escherichia coli] >gpI32495[IECU32495]_1 DegQ [Escherichia coli]		58	35
					transmembrane receptor [Bacillus subtilis] >pirSIC5407[8]			
					chemotaxis transducer homolog TLPA - Bacillus subtilis			
79	20	11050	10535	gpIX5445[9]	tral gene product [Escherichia coli] >pirSIS23001 tral protein - Escherichia coli plasmid RP4 >pirSIC36042 tral protein - Escherichia coli		58	35
85	6	2589	2275	gpIM6451[91]	transport protein [Escherichia coli] >pirSID40840 spermidine/putrescine transport protein D - Escherichia coli		58	38
88	5	2328	3020	gpIL476[7]	path gene product [Vibrio harveyi]		58	37

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

100	17	4026	4235	gpIS5325I	orf V1...orf C4 [tomato leaf curl virus TLCV, Australian isolate, Genomic Complete, 6 genes, 2766 nt]. [Unknown.]	58	47
					>pirSIQ1888 AL2 protein - tomato yellow leaf curl virus (strain Australia)		
100	22	7823	7101	gpIM77039I	E.coli MsbB protein gene, complete cds. [Escherichia coli] >pirSB42608 OrfU upstream of msbB - Escherichia coli (fragment)	58	42
101	33	19952	20701	gpIM57692I	membrane protein for maltose transport [Thermoanaerobacterium thermosulfurigenes] >gpIM57692ITTPULSA_4 membrane protein for maltose transport [Thermoanaerobacterium thermosulfurigenes] >pirSS37704 amyD protein - Thermoanaerobacterium thermosulfurigenes	58	33
104	11	4746	5162	gpIU28377I	Escherichia coli K-12 genome; approximately 65 to 68 minutes. [Escherichia coli]	58	49
112	12	7865	8896	gpIM26929I	D-2-hydroxyisocaproate dehydrogenase [Lactobacillus casei] >gplA1493l[A1493l_1 D-2-Hydroxyisopranseure-dehydrogenase [Lactobacillus casei] >pirSIDELBC_D-2-hydroxyisocaproate dehydrogenase (EC 1.1.1.-) - Lactobacillus casei	58	38
114	3	369	1082	gpIU18997I	Escherichia coli K-12 chromosomal region from 67.4 to 76.0 minutes. [Escherichia coli] >gplU15661IECU15661_2_HhoB [Escherichia coli] >gplU32495IECU32495_2_DegS [Escherichia coli]	58	41
122	6	2873	3442	gpIX56678I	dciAE gene product [Bacillus subtilis]	58	
128	4	4092	3661	gpIL44753I	Hypothetical protein (GB:U14003_130) [Haemophilus influenzae] >gplU32696IHU32696_3 hypothetical protein (GB:U14003_130) [Haemophilus influenzae] >gplU00070IHU00070_14 hypothetical protein (GB:U14003_130) [Haemophilus influenzae] >gplU32805IHU32805_3 in	58	30
					methyl accepting chemotaxis homolog [Treponeema denticola]	58	34
134	3	1234	1025	gpIU33210I	aspartokinase II [Bacillus sp.] >pirSIA43946 aspartate kinase (EC 2.7.2.4) II precursor - Bacillus sp. (strain MGA3)	58	39
136	3	1970	1038	gpIM93419I	OrfUU [Escherichia coli]	58	41
141	8	4280	4942	gpIU23163I	OrfUU [Escherichia coli]	58	43

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

142	2	1395	3965	gpiX59543]	M1 subunit of ribonucleotide reductase [Homo sapiens] >gpiX5961_7 HSRR1LS_1 large subunit ribonucleotide reductase [Homo sapiens] >pirSIS16680 ribonucleoside-diphosphate reductase (EC 1.17.4.1) chain M1 - human		58	41
167	2	732	1049	gpiX835981	elongation factor Ts [Thermus aquaticus thermophilus] >pirSIS51095 elongation factor Ts - Thermus aquaticus	58	30	
179	1	3	452	gpiU242531	folylpolyglutamate synthetase [Homo sapiens] type VI collagen alpha-2 subunit preprotein [Gallus gallus]	58	37	
566	4	995	1270	gpiX150411	>gpiX56659 GDCOL6A2G_1 type VI collagen subunit alpha2 [Gallus gallus] >pirSIS04111 collagen alpha 2(V) chain long form precursor - chicken >gpiX56395 GGCOLVIA_1 type VII collagen alpha-2 subunit	58	41	
634	1	12	467	gpiL44753]	hypothetical protein (GB:U14003_130) [Haemophilus influenzae] >gpiU32696 HTU32696_3 hypothetical protein (GB:U14003_130) [Haemophilus influenzae] >pirU00070 HTU00070_14 hypothetical protein (GB:U14003_130) [Haemophilus influenzae] >gpiU32805 HTU32805_3 in	58	32	
1	11	3440	4156	gpiD835361	proline-tRNA ligase [Escherichia coli]	57	38	
1	17	8526	8290	gpiM32474]	Rattus norvegicus carnoembryonic antigen-related protein (CGM1) mRNA, complete cds. [Rattus norvegicus] >pirSIA35364 carnoembryonic antigen-related protein (clone mCGM1) - rat	57	37	
8	46	27159	27719	gpiU259961	DnaJ [Haemophilus ducreyi]	57	44	
15	2	359	862	gpiM351061	Rat heart-derived c-fos-1 proto-oncogene mRNA, complete cds. [Rattus norvegicus]	57	30	
15	4	1688	2896	gpiX159811	E. coli sbcC gene (ORF-45) for SbcC. [Escherichia coli] >pirSIS0349 hypothetical 45K protein (sbcC 5' region) - Escherichia coli	57	31	
16	3	2001	1369	gpiZ150561	murE gene product [Bacillus subtilis] >pirSIB47691 UDP-N- acetylMuramoyl-L-glutamate--2,6-diaminopimelate ligase (EC 6.3.2.13) - Bacillus subtilis	57	34	
16	47	27505	28170	pirSIS19739	integral membrane protein - Rhodobacter capsulatus	57	31	

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

23	5	1554	2603	gpiM6307I	muts gene product [Azotobacter vinelandii] >pirSIA53296 DNA	57	34
26	26	14567	15556	gpiU3969II	mismatch repair protein MutS - Azotobacter vinelandii methylgalactoside permease ATP-binding protein [Mycoplasma genitalium] >pirSB64213 methylgalactoside permease ATP- binding protein homolog - Mycoplasma genitalium (SGC3) >gpiU02149 MGU02149_1 Mycoplasma genitalium random genomic clone sc8a, partial cds.	57	36
26	29	16354	17610	gpiU3969II	hypothetical protein (SP:P32720) [Mycoplasma genitalium] >pirSID64213 hypothetical protein homolog MG121 - Mycoplasma genitalium (SGC3)	57	31
29	23	10614	10126	gpiX64558I	aprF gene product [Pseudomonas aeruginosa] >pirSIS26698 aprF	57	34
40	3	1163	1582	gpiM77039I	E.coli MshB protein gene, complete cds. [Escherichia coli] >pirSB42608 OrfU upstream of msbB - Escherichia coli (fragment)	57	40
42	23	15434	16132	gpiX72695I	RNA polymerase, beta' subunit (prime) [Thermotoga maritima] >pirSIS41467 DNA-directed RNA polymerase (EC 2.7.7.6) beta' chain (prime) - Thermotoga maritima	57	34
44	23	13698	13420	gpiU29134I	transketolase [Xanthobacter flavus] >gpiU29134 XFU29134_1 transketolase [Xanthobacter flavus]	57	42
55	12	12176	10743	gpiU33007I	D9461.18p [Saccharomyces cerevisiae]	57	40
63	34	21861	20098	pirSIR3EC1	ribosomal protein S1 - Escherichia coli >pirSIS29161 ribosomal protein S1 - Escherichia coli (fragment) {SUB 1-20} >gpiX00785IECRPSA01_2 E. coli rpsA operon leader sequence. [Escherichia coli] {SUB 1-21}	57	38
64	3	884	1783	gpiL46284I	hypothetical protein (GB:D26185_99) [Haemophilus influenzae] >gpiU32838 HTU32838_6 hypothetical protein (GB:D26185_99) [Haemophilus influenzae] >gpiU00086 HTU00086_18 hypothetical protein (GB:D26185_99) [Haemophilus influenzae] >gpiU32785 HTU32785_1 methyl	57	40
70	25	15895	16557	gpiX56234I	cysteine-tRNA ligase [Escherichia coli]	57	42
74	16	9350	9976	gpiU12289I	prolipoprotein diacylglycerol transferase [Escherichia coli]	57	42

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				>pirSIA56149 prolipoprotein diacylglyceryl transferase (EC 2.3.1.1) - Escherichia coli		
76	4	3657	3022	gpiU43739I Borrelia burgdorferi fesmid clone 31, complete sequence. [Borrelia burgdorferi] >gpiU76303IBORFTSA_3 Borrelia burgdorferi ftsA gene, 3' end of cds, ftsZ, orf230, smf, hisVU, flgBCE, fltEFGHI, flbABC genes, complete cds. [Borrelia burgdorferi] >gpiX966685IB unknown [Moraxella catarrhalis]	57	28
81	12	9993	11123	gpiU492269I helicase [Staphylococcus aureus] >pirSIS27667 DNA helicase perA - Staphylococcus aureus >pirSIS39923 perA protein - Staphylococcus aureus	57	32
81	21	18989	17592	gpiM63176I L.delbrueckii nifS-like gene (partial). [Lactobacillus delbrueckii] >pirSIS16047 nifS protein homolog - Lactobacillus delbrueckii unknown [Escherichia coli]	57	37
100	21	5706	7010	gpiX61190I H. influenzae predicted coding region HI0926 [Haemophilus influenzae] >gpiU32774IHU32774_6 H. influenzae predicted coding region HI0926 [Haemophilus influenzae] >gpiU00078IHU00078_71 H. influenzae predicted coding region HI0926 [Haemophilus influenzae]	57	41
101	21	11344	12489	gpiU02965I penicillin-binding protein 1A [Pseudomonas aeruginosa]	57	46
101	26	16171	16473	gpiL45564I TagE [Vibrio cholerae] >gpU390681VCTU39068_4 Vibrio cholerae pathogenicity island, partial and complete cds. [Vibrio cholerae]	57	35
102	6	3868	3302	gpiL13867I >pirSIC2569 tagE protein - Vibrio cholerae (strain 0395)	57	34
102	7	4929	4129	gpiU07173I unknown [Bacillus subtilis]	57	39
103	1	182	703	gpiD26185I recN gene product [Escherichia coli]	57	42
104	20	8066	8869	gpiU36840I GCN2 gene product [Saccharomyces cerevisiae] >pirSIOKBYN2	57	33
121	32	17961	18137	gpiM27082I protein kinase GCN2 (EC 2.7.1.-) - yeast [Saccharomyces cerevisiae] >gpiU51030YSYCD9954_16 Protein kinase, phosphorylates the alpha subunit of eIF-2 (Swiss prot accession number P15442) [Saccharo	57	47
121	33	19184	18273	gpiU00006I E. coli chromosomal region from 89.2 to 92.8 minutes.	57	29

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

128	5	4206	4835	gpiU09711	flagella switch protein [Borrelia burgdorferi] burgdorferi]		57	40
141	4	2425	2622	gpiL211951	>gpiL763031BORFTSA_11 flagellar switch protein [Borrelia burgdorferi]		57	33
144	1	772	62	pirSIA6033	serotonin 5-HT7 receptor protein [Homo sapiens]		57	46
			2	pirSIA6033	80K antigen - Lyme disease spirochete burgdorferi) [SUB 381-475]		57	
150	5	3378	4838	gpiU000191	B2235_C2_195 [Mycobacterium leprae]		57	40
183	1	1	555	gpiX517381	figJ gene product (5 AA) [Salmonella typhimurium] >pirSISMEBH1 flagellar hook-associated protein 1 - Salmonella typhimurium >gplM24466ISTYFLGH_4 S.typhimurium flagellar L-ring (flgH), flagellar P-ring (flgI), and flagellar (flgJ) genes, complete cds. [Sa]		57	28
16	7	5049	4462	gplM248901	esterase [Acinetobacter calcoaceticus]		56	41
16	38	22640	22410	gpiL142851	signaling protein [Plasmid pCF10] >pirSIC53309 prgY protein - Enterococcus faecalis plasmid pCF10		56	35
19	14	13098	10726	gplM242781	lig gene product [Escherichia coli]		56	38
21	6	4067	4681	gpiL772461	ypfQ gene product [Bacillus subtilis]		56	42
22	15	9726	10397	gpiX593991	AdgA protein [Rhodobacter capsulatus] >pirSIS15555 adgA protein - Rhodobacter capsulatus		56	39
29	10	4741	3581	gpiU000131	nifS [Mycobacterium leprae]		56	39
33	1	1	366	gpiU396881	hydroxymethylglutaryl-CoA reductase [Mycoplasma genitalium] >pirSIIID64209 hydroxymethylglutaryl-CoA reductase (NADPH) homolog - Mycoplasma genitalium (SGC3)		56	30
41	4	1447	1704	gpiU000061	E. coli chromosomal region from 89.2 to 92.8 minutes. [Escherichia coli]		56	33
42	12	6057	6617	gpm349951	B.subtilis minor sigma-37 factor of RNA polymerase (rpof, sigB), complete cds. [Bacillus subtilis] >pirSIA36131 hypothetical protein V (sigB 5' region) - Bacillus subtilis >gplL35574[BACRSBU_5 sigma-B positive regulator [Bacillus subtilis] [SUB 1-21]		56	32
44	18	10412	11143	gpm5559171	Spirochaeta aurantia anthranilate synthase component I (lspE)		56	34

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

47	16	10115	9225 gpIX816421	orf gene product [Wolinella succinogenes] >pirSIS50154	56	34	
48	1	765	4 gpIU452851	hypothetical protein - Wolinella succinogenes	56	37	
48	8	5573	6028 gpIM176431	specific [16-kDa vacuolar proton pump subunit [Homo sapiens]]	56	25	
			B.subtilis spoIIA locus sporulation genes. [Bacillus subtilis]	>pirSIS55646 stage II sporulation protein AA protein - Bacillus subtilis	56		
49	11	8342	7806 gpID261851	unknown [Bacillus subtilis]	56	36	
51	5	5745	6236 gpIX168171	Klebsiella pneumoniae gyra gene for DNA gyrase subunit A (EC 5.99.1.3). [Klebsiella pneumoniae]	56	38	
61	9	3590	4795 gpID261851	unknown [Bacillus subtilis]	56	40	
63	24	13179	12289 gpIL103281	f270 gene product [Escherichia coli]	56	34	
64	28	11918	12559 gpIM642731	transfer RNA-Met synthetase [Thermus thermophilus] >pirSISYTWMT methionine-tRNA ligase (EC 6.1.1.10) - Thermus aquaticus	56	36	
74	11	6202	7569 gpIX755681	ICFG [unidentified] >pirSIS348034 Icfg protein - Synecchocystis sp. (PCC 6803) >pirSIS38573 ICFG protein - Synechocystis sp. (strain PCC6803)	56	40	
79	21	10889	11098 gpIM594491	polypeptide chain-binding protein [Zea mays] >gpIM594491MZEB70A_1 polypeptide chain-binding protein [Zea mays] >pirSISQ0966 immunoglobulin-binding protein homolog b70 - maize (fragment) [SUB 1-467]	56	34	
85	5	2314	1505 gpIU2556821	Lpp38 [Pasteurella haemolytica]	56	36	
87	5	3287	2613 gpIX33291	uracil phosphoribosyltransferase [Lactococcus lactis]	56	34	
101	13	6208	5834 gpID261851	replicative DNA Helicase [Bacillus subtilis]	56	28	
101	14	7059	6637 gpIU14031	50S ribosomal subunit protein L9 [Escherichia coli] >gpIX040221ECRPSFRI_4 E. coli genes rpsF, rpsR and rplI for ribosomal proteins S6, S18, L9 [Escherichia coli] >pirSIS56428 50S ribosomal chain protein L9 - Escherichia coli	56	30	
104	21	8701	9282 gpIM302971	B.subtilis recombination and sporulation protein (recN, spoIVB) genes, complete cds, arginine hydroximinate resistance (ahcC) gene, 3' end. [Bacillus subtilis] >pirSISB35128 recN homolog - Bacillus	56	35	

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				subtilis		
106	1	3	740	gplI061481 golgin-165 [Homo sapiens] >pirISIJH0820 160K golgi antigen - human (fragment)	56	33
112	16	10511	11668	gplU465421 ScbA [Streptococcus crista]	56	36
115	2	1153	630	gplX723821 R capsulatus nifR3 DNA. [Rhodobacter capsulatus] >pirISJS34980 hypothetical protein (nifR3_5' region) - Rhodobacter capsulatus	56	43
150	2	3212	2412	gplD281181 DB1 [Homo sapiens] >gplD281181HUMDB1_1 DB1 [Homo sapiens] >pirISIA533772 transcription factor DB1 - human	56	39
150	11	10373	9024	gplU152021 Bacillus subtilis comE operon encoding ORF1, ORF2, ORF3 and Reverse-ORF genes, complete cds. [Bacillus subtilis] >pirIS39864 ComE ORF2 - Bacillus subtilis	56	39
161	4	1128	730	gplS512241 phytochrome [Ceratodon purpureus] >pirISJS27396 phytochrome - (Ceratodon purpureus) (SUB 49-539)	56	37
697	1	213	7	gplM601771 enterobactin [Escherichia coli] >gplM601771ECOENTF_1 enterobactin [Escherichia coli] >pirISIYGCECF enterochelin synthetase (EC 6.4.1.1) component F - Escherichia coli >gplM17354ECOENTFA_1 E.coli entF gene encoding serine activating enzyme, 3' end. [Esche	56	36
733	1	455	243	gplX143091 Murine mRNA for 4F2 antigen heavy chain. [Mus musculus] >pirISJS03600 cell surface antigen 4F2 heavy chain - mouse	56	40
1	18	9516	8524	gplX969831 hypothetical protein [Bacillus subtilis]	55	40
4	2	2681	1215	gplZ681951 Gcd6p [Saccharomyces cerevisiae] >gplL07115YSGCDA_2 guanine nucleotide exchange factor, eIF-2B, delta subunit [Saccharomyces cerevisiae] >gplZ68194ISCS142A_12 Gcd6p [Saccharomyces cerevisiae] >pirISIA48156 translation regulator GCD6 - yeast (Saccharomy	55	27
4	9	6375	7040	gplX871131 DnaJ protein [Agrobacterium tumefaciens]	55	47
5	4	2097	2438	gplL454451 H. influenzae predicted coding region H10807 [Haemophilus influenzae] >gplU32763HTU32763_2 H. influenzae predicted coding region H10807 [Haemophilus influenzae] >gplU00077HTU00077_66 H. influenzae predicted coding region	55	36

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

8	15	7544	8125	gplL45104	HI0807 [Haemophilus influenzae] oxygen-independent coproporphyrinogen III oxidase [Haemophilus influenzae] >gplU32729 HU32729_6 oxygen- independent coproporphyrinogen III oxidase [Haemophilus influenzae] >gplU00073 HU00073_85 oxygen-independent coproporphyrinogen III oxidase [Haemophilus influenzae] >pirSH64222 hypothetical protein MG207 -	55	32
16	10	7078	6437	gplU39698I	M. genitalium predicted coding region MG207 [Mycoplasma genitalium] >pirSH64222 hypothetical protein MG207 -	55	28
16	54	33148	32180	gplX84019I	orf3 gene product [Zymomonas mobilis] >gplX84019 ZMDNAGRP_3 orf3 gene product [Zymomonas mobilis]	55	24
19	4	2825	3775	gplF1453I	histidyl tRNA synthetase [Sus scrofa] >gplF1453I	55	39
21	15	9866	10165	gplL17342I	VaLy-tRNA synthetase [Haemophilus parainfluenzae] >gplL17342I	55	34
22	1	105	1673	gplZ46812I	ZK675_1 [Caenorhabditis elegans] >gplZ46812 CEZK675_1 ZK675_1 [Caenorhabditis elegans]	55	31
27	6	2720	3091	gplJ04243I	unknown protein [Salmonella typhimurium] >pirSIC37890 hypothetical protein (prfA 3' region) - Salmonella typhimurium (fragment)	55	27
28	1	808	137	gplX95575I	cytochrome oxidase subunit 2 [Mitochondrion Chorthippus parallelus] >gplX95575I	55	44
44	19	11435	12082	gplL45882I	H. influenzae predicted coding region HI1248 [Haemophilus influenzae] >gplU32805 HU32805_1 H. influenzae predicted coding region HI1248 [Haemophilus influenzae] >gplU00082 HU00082_17 H. influenzae predicted coding region HI1248 [Haemophilus influenzae]	55	38
44	21	12478	13056	gplU39722I	M. genitalium predicted coding region MG372 [Mycoplasma genitalium] >pirSB64241 hypothetical protein MG372 - Mycoplasma genitalium (SGC3) recombination protein [Bacillus subtilis] >gplX02369 BSORIC_5 Bacillus subtilis oriC region. [Bacillus subtilis] [SUB 48-370]	55	31
51	1	2009	3241	gplD26185I	Bacillus subtilis oriC region. [Bacillus subtilis] [SUB 48-370]	55	31
53	7	3564	3352	gplA7709I	yptA gene product [Bacillus subtilis]	55	27

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

54	25	7958	8197	gp V01146	Genome of bacteriophage T7. [Bacteriophage T7] >pirSISUVBPPB7	55	41
66	1	979	2	gp Z48001	DNA maturase B - phage T7 >pirSIS42338 gene 19 protein -		
70	22	15248	15571	pirSISF53402	orf3 gene product [Thermus aquaticus thermophilus] >pirSIS52278 hypothetical protein 3 - Thermus aquaticus	55	41
80	4	1771	1226	gp U09005	cysteine--tRNA ligase (EC 6.1.1.16) - Bacillus stearothermophilus	55	31
86	1	232	468	gp Z21970	HfIC [Vibrio parahaemolyticus] fragment	55	39
101	31	19067	18792	gp L23195	54CP [Chloroplast Arabidopsis thaliana] >pirSIS36637 signal recognition particle 54CP protein precursor - Arabidopsis thaliana	55	40
					cytoplasmic dynein heavy chain [Drosophila melanogaster] >pirSIA54794 dynein heavy chain, cytoplasmic - fruit fly	55	40
113	1	1	1668	gp X68309	[Drosophila melanogaster] >gp L25122 DRODYNE1H_1 dynein heavy chain [Drosophila melanogaster] [SUB 1877-1998]	55	44
122	5	2451	2308	gp M64780	ERCC3 gene product [Drosophila melanogaster] >pirSIS26719	55	37
135	6	3406	34068	gp L45521	ERCC3 protein - fruit fly [Drosophila melanogaster] agrin [Rattus norvegicus]	55	44
					D-alanine permease [Haemophilus influenzae] >gp U32770 HIU32770_3 D-alanine permease [Haemophilus	55	35
					influenzae] >gp U00078 HIU00078_28 D-alanine permease		
					[Haemophilus influenzae] >gp U32716 HIU32716_5 alanine		
					permease [Haemophilus influenzae] >pirSISH64099		
175	4	2312	1494	gp D26134	297 amino acids peptide, unknown function [Pseudomonas	55	44
183	2	699	917	gp 019171	aeruginosa]		
					52.55kD protein [Human adenovirus type 2] >gp J019171 ADRCG_13 52.5kD protein [Human adenovirus	55	38
					type 2] >pirSISWMA52 late L1 52K protein - human adenovirus 2		
					>gp M73260 ADRCOMPGEN_1 Mastadenovirus h5 gene, complete genome. [Mastadenovirus h5] [SUB 173-4]		
563	1	1	501	gp M24537	GTP-binding protein [Bacillus subtilis] >pirSIS32804 GTP- binding protein, spoOB 3'-region - Bacillus subtilis	55	37
					>gp K02666IBACSP00B2_2 Bacillus subtilis spoOB early		

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				sporulation gene, complete cds. [Bacillus subtilis] [SUB 1-65]		
647	1	582	358	gpIU293991 major outer sheath protein [Treponema denticola]	55	44
1	16	8237	7461	gpIM770391 E.coli MsbbB protein gene, complete cds. [Escherichia coli] >pirSB42608 OrfU upstream of msbB - Escherichia coli (fragment)	54	40
15	6	3330	3791	gpIX142981 Human mRNA for dystrophin. [Homo sapiens] >gplM86903HTMDYST20_1 dystrophin gene product [Homo sapiens] [SUB 2980-3685] >gplL05649HTMDYST15_1 dystrophin [Homo sapiens] [SUB 2850-2979] >pirSIS02109 dystrophin - human (fragment) [SUB 2305-2366] >gplM2326 UDP-MurNac-tripeptide synthetase [Haemophilus influenzae] >gplU32793HTU32793_5 UDP-MurNac-tripeptide synthetase [Haemophilus influenzae] >gplU0008HTU0008_93 UDP-MurNac-tripeptide synthetase [Haemophilus influenzae]	54	26
16	2	1371	1093	gpIL457691 UDP-MurNac-tripeptide synthetase [Haemophilus influenzae] >gplU32793HTU32793_6 UDP-MurNac-tri methyl accepting chemotaxis homolog [Treponema denticola] M.xanthus frzG and frzF genes, complete cds. [Myxococcus xanthus] >pirSIXYZFG frzG protein - Myxococcus xanthus proflaggrin [Homo sapiens]	54	33
19	5	5650	3800	gpIU332101 M.xanthus frzG and frzF genes, complete cds. [Myxococcus xanthus] >pirSIXYZFG frzG protein - Myxococcus xanthus	54	30
19	22	17033	16728	gpIM1352001 Escherichia coli K-12 genome; approximately 65 to 68 minutes. [Escherichia coli]	54	29
26	12	7612	7229	gpIM605031 preprotein translocase secY subunit [Mycoplasma genitalium] >pirSH64218 preprotein translocase secY - Mycoplasma genitalium (SGC3)	54	40
36	21	16662	17306	gpIU283771 M. genitalium predicted coding region MG372 [Mycoplasma genitalium] >pirSB64241 hypothetical protein MG372 - Mycoplasma genitalium (SGC3)	54	20
37	8	2912	3670	gpIU396951 esterase [Bacillus stearothermophilus] >pirSJC1374 carboxylesterase (EC 3.1.1.) - Bacillus stearothermophilus (strain IFO 12550)	54	31
44	22	13047	13349	gpIU397221 S.acidocaldarius ATP synthase (non-FoF1) membrane subunit P (atpP) gene, complete cds. [Sulfolobus acidocaldarius]	54	35
47	29	15854	16273	gpID126811	54	38
50	1	1	231	gpIJ047401	54	33

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				>pirSIA33351 H+-transporting ATP synthase [EC 3.6.1.34]	
63	20	8698	7433	gpi[U00039] proteolipid chain - Sulfolobus acidocaldarius	
				xylose kinase [Escherichia coli] >gpi[K01996]ECOXYLABA_2	
				xyB gene product [Escherichia coli] >gpi[X0469]IECXYLK_2 E.	
				coli genes for xylose isomerase and xylose kinase [Escherichia	
				coli] >pirSIECXY xyulokinase (EC 2.7.1.17) - Escherichia	
				coli	
71	2	1292	567	gpi[U35369] D,D-carboxypeptidase [Enterococcus faecalis]	
100	37	15457	17007	gpi[X91047] hook-associated protein 2 [Xenorhabdus nematophilus]	
101	34	20826	21446	gpi[X66092] C.perfringens ORF for putative membrane transport protein.	
				>pirSIA56641 probable membrane	
				[Clostridium perfringens] >pirSIA56641 probable membrane	
				transport protein - Clostridium perfringens	
112	2	1870	914	gpi[D64006] hypothetical protein [Synechocystis sp.]	
117	6	4154	4732	gpi[M68971] hexokinase type II [Rattus norvegicus] >pirSIS15385 hexokinase	
				(EC 2.7.1.1) II precursor - rat >gpi[M68972]RATHKINAH_1	
				hexokinase type II [Rattus norvegicus] {SUB 402-917}	
				>gpi[Z46367]HSHKEX16_1 hexokinase II [Homo sapiens] {SUB	
				741-791} >pirSIS52114 typ	
121	7	2265	3446	gpi[Z235865] SpoVD [Bacillus subtilis] >pirSIS43863 SpoVD protein - Bacillus	
				subtilis {SUB 1-589} >gpi[L09703]BACPBPSPOV_4 penicillin-	
				binding protein [Bacillus subtilis] {SUB 1-69}	
				>gpi[Z15056]BSSPOG_1 SpoVD [Bacillus subtilis] {SUB 595-645}	
126	2	257	1354	gpi[X96683] phospho-N-acetyl muramoyl-pentapeptide-transferase [Borrelia	
				burgdorferi] >gpi[X96432]BBMRAYFTS_1 phospho-N-	
				acetyl muramoyl-pentapeptide-transferase [Borrelia burgdorferi]	
136	4	1964	2590	gpi[L09228] Bacillus subtilis spoVA to sera region. [Bacillus subtilis]	
				>pirSIS45556 hypothetical protein X14 - Bacillus subtilis	
154	1	887	3	gpi[X75439] isoleucyl tRNA synthetase [Staphylococcus aureus]	
566	1	329	751	gpi[X69292] smooth muscle myosin heavy chain [Homo sapiens]	
677	1	3	506	gpi[U29399] major outer sheath protein [Treponema denticola]	
1	12	4378	5430	gpi[LL12684] Pseudomonas putida recA protein gene, complete cds.	

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

			[Pseudomonas putida]
12	4	3871	2531 gpiX731241 ipa-67d gene product [Bacillus subtilis] >pirISI39722 hypothetical protein - Bacillus subtilis
13	1	666	4 gpiX569581 ankyrin (brank-2) [Homo sapiens]
13	11	6925	6701 gpiU172461 Belongs to the ATP-binding transport protein family (ABC transporters) [Saccharomyces cerevisiae] >pirISI51433 MDL1
16	19	12790	12350 gpiL221671 LZIP-1 and LZIP-2 [Mus musculus]
26	6	4160	4489 gpiX790751 off 208 gene product [Coxiella burnetii] >pirISI34297 hypothetical protein 208 - Coxiella burnetii
32	4	3007	3594 gpiU356731 Borrelia burgdorferi OrfR gene, partial cds, and S20, HBBw, OrfH and Rho genes, complete cds [Borrelia burgdorferi]
35	28	17196	16795 gpiM811681 rubredoxin oxidoreductase [Desulfovibrio vulgaris] >gpiM28848 IDVURUBRB0_1 rubredoxin oxidoreductase [Desulfovibrio vulgaris] >pirISRDDVBX rubredoxin--NAD+ reductase (EC 1.18.1.1) - Desulfovibrio vulgaris
36	19	13115	13906 gpiM2970II S.typhimurium polymerase III polymerase subunit gene, complete cds. [Salmonella typhimurium] >pirISIA45915 DNA-directed DNA polymerase (EC 2.7.7.7) III alpha chain - Salmonella typhimurium >gpiD49445 ECODNAE_1 DnaE, DNA polymerase III holoenzyme catalytic
42	6	3076	36669 gpiL449831 primosomal protein replication factor [Haemophilus influenzae] >gpiU32718HTU32718_10 primosomal protein replication factor [Haemophilus influenzae] >gpiU000721HTU00072_68 primosomal protein replication factor [Haemophilus influenzae] >gpiU32827HTU32827_-
42	10	4662	5078 gpiX024991 Rhodospirillum rubrum atp operon. [Rhodospirillum rubrum] >pirISI08579 hypothetical protein 2 - Rhodospirillum rubrum >gpiX024991RRATP_4 Rhodospirillum rubrum atp operon. [Rhodospirillum rubrum] [SUB 592-811]
43	3	1637	2905 gpiX724911 Ig kappa light chain (V) [Homo sapiens] >pirISI340381 Ig kappa chain - human

TABLE 2.
Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

44	25	14821	14375 gpiU090013	lkt [Mycobacterium leprae]	53	40
54	29	10385	8799 gpiU36427	glutamate synthase small subunit gltD [Thiobacillus ferrooxidans]	53	38
55	14	13304	12729 gpiU15609	flagellar switch protein [Treponema dentitcola] >gpiL3685 lTRPFLIG_1 flag gene product [Treponema dentitcola]	53	27
59	1	1597	359 gpiM77039	E.coli Msbb protein gene, complete cds. [Escherichia coli] >pirlSICB42608 OrfU upstream of msbB - Escherichia coli (fragment)	53	38
70	14	10256	9798 gpiL26916	Pseudomonas aeruginosa (pN) gene, complete cds; ORF1, complete cds; ORF2, complete cds. [Pseudomonas aeruginosa] >pirlSIC53373 hypothetical protein 2 (pN 3' region) - Pseudomonas aeruginosa	53	28
70	18	13714	12437 gpiU21853	unknown [Anabaena sp.]	53	40
74	8	1548	3560 gpiD10280	myosin heavy chain [Oryctolagus sp.]>pirlSIA38650 myosin heavy chain, embryonic smooth muscle - rabbit (fragment)	53	34
77	1	719	3 gpiU39688	hydroxymethylglutaryl-CoA reductase [Mycoplasma genitalium] >pirlSDD64209 hydroxymethylglutaryl-CoA reductase (NADPH) homolog - Mycoplasma genitalium (SGC3)	53	33
79	14	7556	8062 pirlSIC7CHA	spectrin alpha chain, brain - chicken >gpiX14518 GGSPECA5_1 Chicken gene for spectrin alpha-chain 5' end. [Gallus gallus] (SUB 1-27)	53	20
81	11	8080	9603 gpiU49269	amidase [Moraxella catarrhalis]	53	41
81	13	12604	11219 gpiD26185	unknown [Bacillus subtilis]>gpiU02604 BSU02604_3 ORFY [Bacillus subtilis] (SUB 1-260)	53	36
88	6	2957	3565 gpiU18539	FhY [Escherichia coli]	53	25
89	11	9584	9054 gpiX52905	Escherichia coli betT, betI, betB and betA genes. [Escherichia coli] >pirlSIS15179 betT protein - Escherichia coli	53	30
100	39	17891	18439 gpiL19346	Escherichia coli N-acetyl muramoyl-L-alanine amidase (amidB) gene, complete cds. DNA repair protein (mutL) gene, partial cds, and two unidentified cds's. [Escherichia coli] >gpiU14003 ECOUW93_80 yjeE gene product [Escherichia coli] >pirlSIS56393 hypothetical	53	39
100	41	18333	19052 gpiL45029	H. influenzae predicted coding region H10388 [Haemophilus	53	29

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

					[influenzae] >gpiU32722 HTU32722_10 H. influenzae predicted coding region HT0388 [Haemophilus influenzae]
					>gpiU00073 HTU00073_9 H. influenzae predicted coding region HI0388 [Haemophilus influenzae]
121	6	1690	2394	gpiZZ5865	SpoVD [Bacillus subtilis] >pirIS43863 SpoVD protein - Bacillus subtilis [SUB 1-589] >gpiL09703 BACPBPSPVOV_4 penicillin-binding protein [Bacillus subtilis] [SUB 1-69]
					>gpiZ15056BSSPOG_1 SpoVD [Bacillus subtilis] [SUB 595-645]
125	1	1385	2176	gpiD50624	Adenosine Deaminase [Streptomyces virginiae] >gpiD50624 STMVBRA1_1 adenosine deaminase [Streptomyces virginiae] [SUB 1-339]
128	3	4162	2798	gpiX7314	hemolysin [Serpulina hydrotsentiae]
143	1	716	3	gpiI12722	transcription factor IIIB 13 RDa subunit [Saccharomyces cerevisiae] >pirIS1A47453 transcription factor III C chain TFC4 - yeast (Saccharomyces cerevisiae)
143	4	1066	1332	gpiX35548	chromogranin B [Bos taurus]
150	1	2293	2730	gpiU35567	Borrelia burgdorferi OrfR gene, partial cds, and S20, HBbu, OrlfH and Rho genes, complete cds. [Borrelia burgdorferi]
150	3	2709	3380	gpiX07693	protein p67 [Schizosaccharomyces pombe] >pirIS1A30185 nuclear protein nuc2+ - fusion yeast (Schizosaccharomyces pombe)
156	2	58	795	gpiD4591	hypothetical protein [Bacillus subtilis]
358	1	366	617	gpiX82209	MN1 gene product [Homo sapiens]
670	1	327	151	gpiU29399	major outer sheath protein [Treponema dentitcola]
2	14	11472	11702	gpiU19615	LET 858 [Caenorhabditis elegans]
6	4	3991	3425	gpiX77091	trkA gene product [Escherichia coli] >gpiU18997 ECOUW67_214 TrkA protein of the constitutive K+ transport system Trk [Escherichia coli] >gpiX52114 ECTRKAG_3 TrkA protein of the constitutive K+ -transport system Trk [Escherichia coli]
10	9	6165	5398	gpiL47539	cyclic AMP receptor protein [Haemophilus somnus]
11	14	6750	6917	gpiM75136	Ictalurid herpesvirus 1 (channel catfish virus [CCV]), strain

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

11	29	16664	17380	gplL45923	auburn 1, complete genome. [Ictalurid herpesvirus 1] >gpiM75136[HI]ICG_64 Ictalurid herpesvirus 1 (channel catfish virus (CCV)), strain auburn 1, complete genome. [Ictalurid herpesvirus 1] >pif1	52	36
16	52	32135	30837	gplL461861	hypothetical protein (SP:P09171) [Haemophilus influenzae] >gpiU32809[HI]U32809_2 hypothetical protein (SP:P09171) [Haemophilus influenzae] >gpiU00082[HI]U00082_56 hypothetical protein (SP:P09171) [Haemophilus influenzae] >gpiU32754[HI]U32754_5 pseudoU synthetase	52	21
19	7	7636	5668	gpiU33210	H. influenzae predicted coding region HI1555 [Haemophilus influenzae] >gpiU32830[HI]U32830_11 H. influenzae predicted coding region HI1555 [Haemophilus influenzae] >gpiU00085[HI]U00085_26 H. influenzae predicted coding region HI1555 [Haemophilus influenzae]	52	21
21	2	1556	1332	gpiX744681	methyl accepting chemotaxis homolog [Treponema denticola] late protein [Human papillomavirus type 15] >pirfSIS36478 late protein - human papillomavirus type 15 >gpiM96284[PPHLDF_1_L] gene product [Human papillomavirus type 15] {SUB 312-355} sucrose synthase [Arabidopsis thaliana]	52	31
23	1	205	2	gpiX70901	dreb1n E2 [Homo sapiens] >gpiID17530[HUMDRE_1_dreb1n E2] [Homo sapiens] >pirfSIN0869 dreb1n E (Clone gDbh13) - human membrane associated ATPase [Haemophilus influenzae] >gpiU32835[HI]U32835_5 membrane associated ATPase [Haemophilus influenzae] >gpiU00085[HI]U00085_85 membrane associated ATPase [Haemophilus influenzae] >pirfSIA6133	52	41
30	11	6103	6936	gpiU008021	dreb1n E2 [Homo sapiens] >gpiID17530[HUMDRE_1_dreb1n E2] [Homo sapiens] >pirfSIN0869 dreb1n E (Clone gDbh13) - human membrane associated ATPase [Haemophilus influenzae] >gpiU32835[HI]U32835_5 membrane associated ATPase [Haemophilus influenzae] >gpiU00085[HI]U00085_85 membrane associated ATPase [Haemophilus influenzae] >pirfSIA6133	52	25
42	3	1276	950	gplL462481	NAD-dependent methylenetetrahydrofolate dehydrogenase-methylenetetrahydrofolate cyclohydrolase [Drosophila melanogaster] >gpiSS59910[S59910_1 NAD-dependent methylenetetrahydrofolate dehydrogenase [Drosophila melanogaster] >pirfSIS32562 methylenetetrahydrofolate	52	40
54	20	7489	6683	gplL079581	orf4 [Bacillus subtilis] >gpiD37799[BACAMOKOO_6 orf4] [Bacillus subtilis]	52	33
55	9	8295	8795	gpiD377991		52	30

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

35	13	13002	12193	gpIU09711	flagella switch protein [Borrelia burgdorferi] >gplL76303IBORFTSA_11 flagellar switch protein [Borrelia burgdorferi]	52	29
63	31	18372	17776	gpIX908721	gp2512 gene product [Homo sapiens]	52	35
68	7	2743	3408	pirSIJN0695	tributyltin chloride resistant protein - Alteromonas sp. (strain M-1)	52	35
70	20	14776	13976	gpIM318271	Bacillus subtilis (clone lambda-BSt) cell division and sporulation protein (dds) gene, complete cds. [Bacillus subtilis]>pirSIJA43727 probable division initiation regulatory protein 1 - Bacillus subtilis >epIM31800IBACDIV_1 B.subtilis division initiation	52	42
77	11	5580	5026	gpIL382521	Acinetobacter lwoffii orf1 and esterase (est) genes, complete cds. [Acinetobacter lwoffii]	52	34
83	11	3711	4682	gpIZ543281	unknown [Schizosaccharomyces pombe]	52	34
104	5	1730	1915	gpIS75651	PSI [Drosophila]>gpIS75661S75666_1 PSI [alternatively spliced] [Drosophila, Pre-mRNA, 3961 nt]. [Drosophila]	52	41
104	13	5305	5478	gpIU000241	u00024q [Mycobacterium tuberculosis]	52	39
112	14	9864	9097	gpIM261301	S. parasanguis adhesin (fimA), ORF1, ORF3, and ORF5 genes, complete cds. [Streptococcus parasanguis]	52	27
115	4	1919	1473	gpIZ562801	carD Gene product [Myxococcus xanthus]	52	35
170	4	3357	2311	gpIU437391	FtsA [Borrelia burgdorferi]>gpIX96433BBFTSWQA_3 ftsA gene product [Borrelia burgdorferi] (SUB 1-36)	52	32
182	2	1189	215	pirSIIB3687	orf34' 5' of tgs - Escherichia coli	52	38
619	1	248	24	gpIZ173721	M.smeognatis asd, ask-alpha, and ask-beta genes. [Mycobacterium smeognatis]>pirSI31804 hypothetical protein y - Mycobacterium smeognatis	52	27
631	1	403	621	gpID640031	hypothetical protein [Synechocystis sp.] >gplD640031SYCSLLE_50 hypothetical protein [Synechocystis sp.]	52	30
15	7	4152	5183	gpIL139731	Gallus gallus EDT-soluble/130 kDa protein mRNA, complete cds. [Gallus gallus]>pirSI47168 cardiac morphogenesis protein ES/130 - chicken (fragment)	51	26

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

17	1	552	73	gpiU3656I	fus-like protein [Homo sapiens]	51	42
19	2	2349	556	gpiU00022I	u0308b [Mycobacterium leprae]	51	37
30	12	7372	8253	gpiX86903I	peptidyl prolyl isomerase [Triticum aestivum] >pirIS55383	51	29
30	14	9302	10168	gpiM31792I	E.coli MreB protein gene, 3' end, MreC protein gene, complete cds, and MreD protein gene, complete cds. [Escherichia coli] >gpiU18997IECOUW67_180 mreC gene product [Escherichia coli] >pirISUV0059 mreC protein - Escherichia coli >gpiM22055IECOMREB_3_E.col	51	39
34	2	623	1327	gpiX16335I	Klebsiella pneumoniae rpoN gene 3'downstream region. [Klebsiella pneumoniae] >pirIS07661 hypothetical protein 162 (rpoN 3' region) - Klebsiella pneumoniae	51	32
37	11	6299	4824	pirSIA2294_0	keratin, 67K type II cytoskeletal - human	51	40
39	1	188	3	gpiS78086I	protein-tyrosine phosphatase [Homo sapiens] >pirSIB44929 protein-tyrosine-phosphatases (EC 3.1.3.48) BPTP-2 - human (fragment)	51	43
39	5	2837	1236	gpiZ22177I	F54G8.4 [Caenorhabditis elegans]	51	35
42	5	2265	3104	gpiM84415I	DNA polymerase [Bacteriophage SP01] >pirSIC1269 DNA-directed DNA polymerase (EC 2.7.7.) - phage SP01	51	30
53	8	6786	7085	gpiL03292I	taurine/beta-alanine transporter [Mus cookii] >pirSIA47194 taurine and beta-alanine transporter. TAUT - mouse	51	37
56	7	3620	3847	gpiK03277I	Tm2 gene product [Drosophila melanogaster] >pirSIA25624 tropomyosin I, embryonic - fruit fly (Drosophila melanogaster)	51	25
66	4	3219	1651	gpiX78998I	endosomal protein [Homo sapiens] >pirIS44243 endosomal protein - human	51	22
75	6	2330	3406	gpiX73124I	ipa-83d gene product [Bacillus subtilis] >pirIS39738 hypothetical protein - Bacillus subtilis	51	34
89	4	1689	2285	gpiU32702I	DNA helicase [Haemophilus influenzae]	51	39
89	12	11391	9544	gpiL46177I	lic-I operon protein [Haemophilus influenzae] >gpiU32829HTU132829_8 lic-1 operon protein [Haemophilus influenzae] >gpiU00085HTU00085_10 lic-1 operon protein	51	31

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				[Haemophilus influenzae] >gpiU32775HU32775_8		
				nucleotidyltransferase [Haemophilus influenzae] >p		
91	4	1983	1642	gplJ15626I dynein [Saccharomyces cerevisiae]	51	25
98	1	754	2	gplD64006I hypothetical protein [Synechocystis sp.]	51	27
101	2	416	1012	gplL45065I hypothetical protein (SP:P33635) [Haemophilus influenzae] >gpiU32726HU32726_5 hypothetical protein (SP:P33635) [Haemophilus influenzae] >gpiU00073HU00073_46 hypothetical protein (SP:P33635) [Haemophilus influenzae] >gpiU32834HU32834_10 rRNA methylas	51	30
101	17	8473	8165	gplJ39923I 30S ribosomal protein S6 [Mycobacterium leprae]	51	30
101	35	21692	22174	gplU10427I CapD [Staphylococcus aureus]	51	28
102	3	2448	1180	gplD64000I hypothetical protein [Synechocystis sp.]	51	38
117	4	2669	2142	gplL45352I trigger factor [Haemophilus influenzae] >gpiU32754HU32754_4	51	27
				trigger factor [Haemophilus influenzae] >gpiU00076HU00076_52 trigger factor [Haemophilus influenzae] >gpiU32700HU32700_4 peptidyl-prolyl cis-trans isomerase [Haemophilus influenzae] >pirISI		
128	2	2982	1534	gplL44753I hypothetical protein (GB:U14003_130) [Haemophilus influenzae] >gpiU32696HU32696_3 hypothetical protein (GB:U14003_130) [Haemophilus influenzae] >gpiU00070HU00070_14 hypothetical protein (GB:U14003_130) [Haemophilus influenzae] >gpiU32805HU32805_3 in	51	26
139	2	851	288	gplD63999I hypothetical protein [Synechocystis sp.]	51	35
572	1	519	235	gplL15202I Bacillus subtilis comE operon encoding ORF1, ORF2, ORF3 and Reverse-ORF genes, complete cds. [Bacillus subtilis] >pirIS39864 ComE ORF2 - Bacillus subtilis	51	33
8	3	3906	2470	gplU23181I ZK84.1 gene product [Caenorhabditis elegans]	50	28
11	11	5426	5923	gplJ45290I rep helicase, single-stranded DNA-dependent ATPase [Haemophilus influenzae] >gpiU32748HU32748_3 rep helicase, single-stranded DNA-dependent ATPase [Haemophilus influenzae] >pirSID64084 rep helicase, single-stranded DNA-dependent ATPase (rep) homolog -	50	32

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

11	18	8259	9131	[gpiU293991] major outer sheath protein [Treponema denticola]			50	45
17	5	2963	2391	[gpiU000081] yejD [Escherichia coli]			50	33
22	14	9154	9777	[gpiX593991] AdgA protein [Rhodobacter capsulatus]	>pirSIS15555 adgA		50	35
22	22	14326	14496	[gpiX023071] E. coli aspA gene for aspartase (L-aspartate ammonia-lyase) (EC 4.3.1.1). [Escherichia coli]			50	37
24	1	1	705	[gpiM1652] growth factor [Drosophila melanogaster]			50	33
27	8	4573	3698	[gpiL462901] ribonucleoside diphosphate reductase B2 subunit [Haemophilus influenzae] >gpiU32839[HTU32839_2 ribonucleoside diphosphate reductase B2 subunit [Haemophilus influenzae]]	>gpiU00086[HTU00086_23 ribonucleoside diphosphate reductase B2 subunit [Haemophilus influenzae]]		50	37
39	10	4905	6335	[gpiU113845] acr206 gene product [Rhizobium meliloti]	>pirSIS49802 hypothetical		50	24
43	5	4605	3028	[gpiZ467291] unknown [Saccharomyces cerevisiae] >pirSIS49802 hypothetical	protein YM9958_03c - yeast (Saccharomyces cerevisiae)		50	31
45	18	15712	16710	[gpiU0000101] transport protein (similarity to antibiotic transport protein acil1_3 from S.coelicolor) [Mycobacterium leprae]			50	24
53	10	5460	5296	[gpiX78057] carletilulin [Zea mays]			50	31
59	5	2553	2170	[gpiU283751] single-stranded DNA-specific exonuclease [Escherichia coli]			50	28
60	1	1	189	[gpiU23764] TonB [Pseudomonas aeruginosa]			50	32
61	6	1735	2823	[gpiL44751] hypothetical protein (GB:L01112_7) [Haemophilus influenzae]	>gpiU32696[HTU32696_1 hypothetical protein (GB:L01112_7)]		50	32
				[Haemophilus influenzae] >gpiU00070[HTU00070_12 hypothetical protein (GB:L01112_7) [Haemophilus influenzae]]	>gpiU32805[HTU32805_1 H. influ			
62	2	954	403	[gpiD64001] hypothetical protein [Synechocystis sp.]			50	37
63	9	3455	3814	[pirSIS53457] dominant autoantigen gp 330 - rat (fragment)			50	50
70	32	19744	20682	[gpiD640061] hypothetical protein [Synechocystis sp.]			50	31
71	1	555	97	[gpiX92441] S cerevisiae 33kb fragment from the right arm of chromosome XV. [Saccharomyces cerevisiae]			50	27
72	4	4363	2789	[gpiX66793] sigma factor 54 [Alcaligenes eutrophus] >pirSIB48362 transcription initiation factor sigma 54 - Alcaligenes eutrophus			50	29

TABLE 2.
Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

89	10	9208	8306	gpiX529051	[Escherichia coli betT, betI, betB and betA genes. [Escherichia coli]]	50	37
96	2	173	658	gpiL456451	>pirSIS15179 betT protein - Escherichia coli hypothetical protein (GB:D10483_22) [Haemophilus influenzae] >gpiU32781IHU32781_10 hypothetical protein (GB:D10483_22) [Haemophilus influenzae] >pirSIA64164 hypothetical protein HII1007 - Haemophilus influenzae (strain Rd KW20) >gpiU00079IHU00079_67 hyp	50	32
100	9	1781	2587	gpiX603951	Hox 4.6 gene product [Mus musculus] >gpiX71422MMHOXD11_1 HOXD-11 [Mus musculus] (SUB 14-336)	50	47
135	4	1991	2278	pirSIS29717	adenylate cyclase (EC 4.6.1.1) type 5 - rat phenylalanyl-tRNA synthetase alpha subunit [Saccharomyces cerevisiae]	50	34
143	5	5051	3876	gpiX946071	endopeptidase I gene product [Bacillus sphaericus] >pirSIS333310 endopeptidase I - Bacillus sphaericus >gpiX69895BSPROTXA_1 peptidase I gene product [Bacillus sphaericus] (SUB 375-396)	50	27
581	1	20	454	gpiX695071	hypothetical protein (GB:U14003_130) [Haemophilus influenzae] >gpiU32696IHU32696_3 hypothetical protein (GB:U14003_130) [Haemophilus influenzae] >gpiU00070IHU00070_14 hypothetical protein (GB:U14003_130) [Haemophilus influenzae] >gpiU32805IHU32805_3 in	50	40
639	1	1	489	gpiL447531	C.burnetii trxB, spoIIIE and serS genes. [Coxiella burnetii] >pirSIS43133 hypothetical protein - Coxiella burnetii >pirSIS31760 hypothetical protein Y - Coxiella burnetii	50	26
2	1	1	708	gpiX756271	folyl-polyglutamate synthetase [Bacillus subtilis] >pirSISB40646 folC - Bacillus subtilis	49	29
2	10	7878	6370	gpiL045201	lysyl-tRNA synthetase analog [Escherichia coli] >pirSIS56383 lysyl-tRNA synthetase genX - Escherichia coli >gpiX59988IECGENXLTR_1 genX gene product [Escherichia coli] (SUB 11-335)	49	32
21	3	1389	2453	gpiUT40031	nifU [Mycobacterium leprae]	49	27
29	8	3666	3211	gpiU000131	nifU [Mycobacterium leprae]	49	30
30	15	10668	12164	gpiL446761	penicillin-binding protein 2 [Haemophilus influenzae]	49	34

TABLE 2.
Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				>gp U32688 HTU32688_11 penicillin-binding protein 2 [Haemophilus influenzae] >gp U00069 HTU00069_30 penicillin-binding protein 2 [Haemophilus influenzae]		
43	1	1	1224	gplM133591 E.coli sppA gene encoding protease IV; complete cds. [Escherichia coli] >pirSIPRECT4_protease IV (EC 3.4.-.-) - Escherichia coli >gp U13772 ECU13772_1 protease IV [Escherichia coli] {SUB 110-433}	49	37
44	13	7058	7838	gplI452061 hypothetical protein (SP:P32662) [Haemophilus influenzae] >gp U32738 HTU32738_4 hypothetical protein (SP:P32662) [Haemophilus influenzae] >gp U00074 HTU00074_85 hypothetical protein (SP:P32662) [Haemophilus influenzae] >gp U32847 HTU32847_4 hydrolase (pho	49	31
47	14	82555	7647	gplM929051 calcium channel alpha-1 subunit [Rattus norvegicus] >pirSIF33610 npD	49	36
60	3	1212	586	gplD174621 Na+-ATPase subunit D [Enterococcus hirae] >pirSIF33610 npD	49	25
79	17	8989	100355	gplL145561 glutamate synthase [Escherichia coli] >gp U000061ECOUW89_26 glutamate synthase [Escherichia coli] >gp V00347 ECRRNB_1E. coli 3' noncoding region preceding the gene rmB which codes for 16S ribosomal RNA. [Escherichia coli] >gp V00348 ECRRNBZ_1E. coli ri	49	34
15	3	855	1724	gplJ039981 P.falciparum glutamatic acid-rich protein 8nen, complete cds. [Plasmodium falciparum] >pirSIA54514 glutamic acid-rich protein precursor - Plasmodium falciparum	48	21
16	6	4339	4063	gplX651951 N-acetylphosphonothricin-tripeptide-deacetylase [Streptomyces viridochromogenes] >pirSIS20686 N-acetylphosphonothricin-tripeptide-deacetylase - Streptomyces viridochromogenes >gplM22827 STMPAT_3 Streptomyces viridochromogenes phosphonothricin N-acetyltran	48	35
16	39	23653	22703	gplU213201 K04G7_3 gene product [Caenorhabditis elegans]	48	27
26	36	23309	24658	gplU13961I protective surface antigen D15 [Haemophilus influenzae] >pirSLJC4078 D-15 protective surface antigen - Haemophilus influenzae	48	29

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

29	25	11464	12030 gplU202171	fibrillin-2 [Mus musculus]	48	36
32	3	1511	3190 gplL449471	apolipoprotein N-acyltransferase [Haemophilus influenzae] >gplU32716 HU32716_3 apolipoprotein N-acyltransferase [Haemophilus influenzae] >gplU00072 HU00072_32	48	30
				apolipoprotein N-acyltransferase [Haemophilus influenzae] >gplU32825 HU32825_3 apolipoprotein		
41	7	1910	2155 gplL368291	alphaA-crystallin-binding protein I [Mus musculus] >gplX68946 MACRYBP1_1 alphaA-CRYBP1 [Mus musculus] (SUB 2024-2688)	48	45
42	16	9061	9687 gplD503031	Ribosomal Protein L10 [Bacillus subtilis]	48	29
47	23	13542	12838 gplL448481	hypothetical protein (SP:P21504) [Haemophilus influenzae] >gplU32705 HTU32705_7 hypothetical protein (SP:P21504) [Haemophilus influenzae] >gplU00071 HTU00071_18 hypothetical protein (SP:P21504) [Haemophilus influenzae] >gplU32814 HTU32814_7 H.influenzae	48	30
49	10	7561	7097 gplD261851	expressed at the end of exponential growth under conditions in which the enzymes of the TCA cycle are repressed [Bacillus subtilis] >gplX16518 BSTMSPRS_3 B.subtilis prs,tms, and ctc (partial) genes for PRPP synthetase and two undefined gene products. [Bacil	48	26
63	6	2357	1895 gplZ331261	membrane forming protein [Mycoplasma capricolum] >pirIS48611 hypothetical protein - Mycoplasma capricolum (SGC3) (fragment) [SUB 1-101]	48	26
85	2	336	764 gplL334681	thermoregulated motility protein [Yersinia enterocolitica type O8] >gplL334681 YEP7TMA_1 thermoregulated motility protein [Yersinia enterocolitica (type O8)]	48	24
100	20	5688	4987 gplS724421	peptidyl-prolyl cis/trans isomerase [Legionella pneumophila] >pirIS30591 outer membrane protein mp precursor - Legionella pneumophila	48	31
101	5	1378	1841 gplM962341	oxaloacetate decarboxylase [Salmonella typhimurium]	48	35
117	9	6598	6206 gplL406321	ankyrin 3 [Mus musculus]	48	30
564	1	2	385 gplL449461	hemolysin [Haemophilus influenzae] >gplU32716 HU32716_2	48	25

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				[hemolysin [Haemophilus influenzae] >gpiU000721HTU00072_31		
				hemolysin [Haemophilus influenzae] >gpiU32825HTU32825_2		
1	26	11879	12952	gplA63161 cyathionine beta-synthase [Haemophilus influenzae] >priSIC64060 hemolysin (tlyC)	47	27
				nitrogen fixation protein [Haemophilus influenzae] >gpiU32841HTU32841_8 nitrogen fixation protein [Haemophilus influenzae] >gpiU000861HTU00086_48 nitrogen fixation protein [Haemophilus influenzae] >gpiU32788HTU32788_4 membrane protein, NADH:ubiquinone o		
8	16	8098	8844	gplA51041 oxygen-independent coproporphyrinogen III oxidase [Haemophilus influenzae] >gpiU32729HTU32729_6 oxygen-independent coproporphyrinogen III oxidase [Haemophilus influenzae] >gpiU00073HTU00073_85 oxygen-independent coproporphyrinogen III oxidase [Haemophilus	47	29
16	37	22283	22083	gpIZ334131 unknown [Pseudomonas syringae] >gpiZ334131PSFOSCG_3 unknown [Pseudomonas syringae] >priSIS44937 hypothetical protein - Pseudomonas syringae	47	23
27	7	3449	3781	gpIZ497821 ywKE gene product [Bacillus subtilis] >priSIS553438_ywKE protein - Bacillus subtilis	47	36
29	13	5615	4911	gpiU000131 dps1 [Mycobacterium leprae]	47	23
37	14	7473	8843	gplA77091 yptA gene product [Bacillus subtilis]	47	24
55	11	10980	9094	gpiU339071 D9461_18p [Saccharomyces cerevisiae]	47	30
57	13	9939	10847	gpIX642591 N-acetylglucosaminyl transferase [Bacillus subtilis] >priSUCI1275 phospho-N-acetylglucosaminyl-pentapeptide-transferase (EC 2.7.8.13) - Bacillus subtilis >gplM31827BACDDSA_1 Bacillus subtilis (clone lambda-BS1) cell division and sporulation protein (dds) ge	47	27
63	8	4537	3023	gplM13169 high affinity ribose transport protein [Escherichia coli] >priSIB26304 ribose transport protein rbsA - Escherichia coli	47	27
104	18	7200	8012	gpiU000211 Mycobacterium leprae cosmid L247. [Mycobacterium leprae]	47	30
158	1	1331	15	gpIU217341 UNC44 [Caenorhabditis elegans] >pirSAS57282 ankyrin-related protein unc-44 - Caenorhabditis elegans (fragment)	47	34

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

175	3	1415	729	gpiU349231	Tbi1 [Escherichia coli]		47	26
326	1	322	2	gpiM831961	microtubule-associated protein IA [Rattus norvegicus] >pirS1A43359 microtubule-associated protein MAP1A - rat		47	38
36	22	17495	18919	gpiX773951	uridine kinase [Saccharomyces cerevisiae] >gpiX539981SCURK1_2 uridine kinase [Saccharomyces cerevisiae] >pirS129374 uridine kinase (EC 2.7.1.48) - yeast (Saccharomyces cerevisiae)		46	27
37	7	2757	3395	gpID308081	Glutathione-regulated potassium efflux system (KcfC) [Bacillus subtilis]		46	30
135	3	3028	878	gpiL459701	hypothetical protein (GB:L12968_1) [Haemophilus influenzae] >gpiL460941HEAH11463_1 hypothetical protein (GB:L12968_1) [Haemophilus influenzae] >gpiU328131HTU32813_2 hypothetical protein (GB:L12968_1) [Haemophilus influenzae] >gpiU328241HTU32824_5 hypothet		46	34
1	15	6998	7420	gpiM164891	Escherichia coli tolQRA gene cluster DNA. [Escherichia coli] >pirISWMEC15_15.5K protein (tolAB operon 5' region) - Escherichia coli >gpiU309341ECU30934_2 Escherichia coli cytochrome oxidase d subunit II (cydB) gene, partial cds, and orf in tolQRA region,		45	30
10	8	5476	5099	gpiZ482391	orf7 gene product [Saccharomyces cerevisiae] >pirS157679 hypothetical protein 7 - yeast (Saccharomyces cerevisiae)		45	24
12	6	4020	4358	gpiU347741	ankyrin-like repeat protein [Orf virus] >gpiS785161S78516_1 G1L gene product [Unknown.]		45	36
17	14	7538	7765	gpiZ482361	X-prolyl dipeptidyl aminopeptidase [Lactobacillus helveticus]		45	27
22	16	10438	11088	gpiL462591	dedA protein [Haemophilus influenzae] >gpiU328361HTU32836_3 dedA protein [Haemophilus influenzae] >gpiU0085HTU00085_96 dedA protein [Haemophilus influenzae] >gpiU327821HTU32782_16 alkaline phosphatase-like protein [Haemophilus influenzae] >pirS1D64133		45	27
44	12	5893	7155	gpiL448831	H. influenzae predicted coding region H10238 [Haemophilus influenzae] >gpiU327101HTU32710_4 H. influenzae predicted coding region H10238 [Haemophilus influenzae] >gpiU00071HTU00071_54 H. influenzae predicted coding region		45	22

TABLE 2.
Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				HI0238 [Haemophilus influenzae]		
81	17	14515	16290	gpiZ232522 [Lactococcus lactis] >pirIS49150	45	28
100	1	1	1059	gpiZ2664991 [Lactococcus lactis]	45	39
104	17	6278	7399	gpiX731241 T01B7.8 [Caenorhabditis elegans]	45	24
13	10	6339	8366	gpiU136751 ipa-65d gene product [Bacillus subtilis] >pirIS339720 hypothetical protein - Bacillus subtilis	45	
26	35	22077	23420	gpiU139611 lactose permease [Citrobacter freundii] >pirSJC2544 lactose carrier protein - Citrobacter freundii	44	31
35	10	6451	7464	gpiD901091 protective surface antigen D15 [Haemophilus influenzae] >pirSJC4078 D-15 protective surface antigen - Haemophilus influenzae	44	29
41	1	23	751	gpiM219941 Rat mRNA for long-chain acyl-CoA synthetase (EC 6.2.1.3). [Rattus norvegicus] >gpiM55642/RATCOAA_1 acyl-CoA synthetase [Rattus norvegicus] >pirSIA36275 long-chain-fatty-acid-CoA ligase (EC 6.2.1.3) - rat	44	30
47	13	7765	6896	gpiS444261 E. coli cysK gene, 3' end, ptsH, ptsI, and crt phototransferase system genes, complete cds. [Escherichia coli] >gpiQ27961/ECOPTSHI_2 ptsI gene product [Escherichia coli] >pirSWQECPII phosphotransferase system enzyme I (EC 2.7.3.9) - Escherichia coli	44	27
98	5	4080	4325	gpiZ115821 cytB gene product [Synechococcus sp. (PCC 7942) (fragment)]	44	32
117	7	4540	5535	gpiZ541421 nuf1 gene product [Saccharomyces cerevisiae] >gpiX73297/SCSETRP4_2 SPC110/NUF1 gene product [Saccharomyces cerevisiae] >gpiU28372/YSYSCD9476_7 Probable essential component of the nucleoskeleton (Swiss Prot accession number P32380) [Saccharomyces cerevisiae]	44	24
152	3	2944	2366	gpiM352001 unknown [Schizosaccharomyces pombe] >gpiX92894/SPHXXIP_1 hexokinase 1 [Schizosaccharomyces pombe]	44	27
				M. xanthus frzG and frzF genes, complete cds. [Myxococcus xanthus] >pirSIXYYZFG frzG protein - Myxococcus xanthus	44	29

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

153	1	748	2 gpxX8274891	pyruvate,orthophosphate dikinase [Mesembryanthemum crystallinum] >gplX824891 CPDD_1 pyruvate,orthophosphate dikinase [Mesembryanthemum crystallinum] >pir SIS55478 pyruvate,orthophosphate dikinase (EC 2.7.9.1) - common ice plant	44	36
156	1	2	340 gplI14321I	BICP4 [Bovine herpesvirus type 1]	44	38
183	3	806	1117 gplU02514I	hook-associated protein 3 [Escherichia coli] >pir S44022 hook-associated protein 3 - Escherichia coli	44	30
1	25	11110	11946 gpxX72888I	mfC gene product [Rhodobacter capsulatus] >pir SIS39893 mfC protein - Rhodobacter capsulatus >gplX79064 RCRNFCD_1 mfC gene product [Rhodobacter capsulatus] [SUB 39-519]	43	28
11	12	5921	6973 gpm63489I	ATP-dependent nuclelease [Bacillus subtilis] >pir SIB39432 ATP-dependent exonuclease synthetase protein AddA - Bacillus subtilis	43	26
42	14	8004	7747 gplU24265I	special lobe-specific protein [Chironomus thummi]	43	33
64	30	12557	13519 gplI45532I	membrane fusion protein [Haemophilus influenzae] >gplU32771 HUU32771_6 membrane fusion protein [Haemophilus influenzae] >gplU32717 HUU32717_6 permease [Haemophilus influenzae] >pir SIG64100 membrane fusion protein (mfC) homolog - Haemophilus influenzae (43	22
77	8	3699	4703 pir SIS027083	tropomin T - fruit fly (Drosophila melanogaster)	43	29
85	9	3407	3718 gpxX64346I	Herpesvirus saimiri complete genome DNA. [Saimiri herpesvirus 1] >pir SIA3681_1 hypothetical protein ORF48 - saimiriine herpesvirus 1 (strain 11) >gplM86409 HSV3PRGEN_1 Herpesvirus saimiri the most three prime end of the genome. [Herpesvirus saimiri] [SU]	43	35
159	1	2	1165 gplI35574I	sigma-B regulator [Bacillus subtilis]	43	20
162	1	302	1714 gplI23195I	cytoplasmic dynein heavy chain [Drosophila melanogaster] >pir SIA54794 dynein heavy chain, cytoplasmic - fruit fly (Drosophila melanogaster) >gplI25122 DRODYNENH_1 dynein heavy chain [Drosophila melanogaster] [SUB 1877-1998]	43	26
1	13	5428	6120 gpxZ67757I	unknown [Schizosaccharomyces pombe]	42	31
6	5	4911	3877 gpxX77091I	trkA gene product [Escherichia coli] >gplU18997 ECOUW67_214	42	28

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

				TrkA protein of the constitutive K ⁺ transport system Trk [Escherichia coli] >gpiIX52114 ECTRKAG_3 TrkA protein of the constitutive K ⁺ -transport system Trk [Escherichia coli]		
>priSS36252 trkA pr						
16	5	2829	4070	gplLA2115 insulin-activated amino acid transporter [Mus musculus] >priSJC4149 adipocyte amino acid transporter - mouse	42	24
100	4	1738	1061	gplIX88798 AT1 gene product [Oryza sativa] >priSS57459 hook-containing protein AT1 - rice	42	37
29	14	5713	5516	gplU0394 envelope protein [Simian immunodeficiency virus] >priSIS46335 envelope protein - simian immunodeficiency virus	41	27
53	8	3943	3542	gplD10474 ORF248 [Synechocystis sp.] >priSJT0603 hypothetical 27.8K protein (frxC 3' region) - Synechocystis sp. (PCC 6803)	41	26
65	3	356	505	gplIA15841 proteinase gene product [Lactococcus lactis cremoris]	41	29
68	2	1774	572	gplIX06545 E. coli genes hsdR and hsdM. [Escherichia coli] >priSINDECKR type I site-specific deoxyribonuclease (EC 3.1.21.3) EcoK chain R - Escherichia coli	41	24
70	30	18836	20197	gplU00001 CDC27 [Homo sapiens]	41	23
101	20	10060	11289	gplL31959 Mus musculus (strain C3H/FL) ORF mRNA, complete cds. [Mus musculus]	41	26
2	6	3975	4871	gplD26185 high level kasugamycin resistance [Bacillus subtilis]	40	27
16	1	993	4	gplU41010 T05A12.2 gene product [Caenorhabditis elegans]	40	27
47	7	4032	2482	gplIX06165 Yeast CDC16 gene. [Saccharomyces cerevisiae] >gplZ280221SCYKL022C_1 CDC16 gene product [Saccharomyces cerevisiae] >priSIA27832 cell division control protein CDC16 - yeast (Saccharomyces cerevisiae)	40	27
11	17	7969	7522	gplY00063 Plasmid falciparum mRNA fragment for knob protein. [Plasmid falciparum] >priSIS14431 Knob-associated histidine-rich protein - Plasmid falciparum (fragment)	39	25
234	1	3	302	gplM87634 BF-1 [Rattus norvegicus] >priSJH0672 brain factor 1 protein - rat	39	32
56	12	5977	6690	gplIX65165 extensin [Volvox carteri] >priSIS22697 extensin - Volvox carteri (fragment)	38	33
70	5	2560	4332	gplD43935 Xanthine Dehydrogenase [Bombyx mori]	36	21

TABLE 2.

Treponema pallidum - Putative coding regions of novel proteins similar to known proteins

56	13	6826	6131	epIU226801	X box repressor [Homo sapiens]	35	19
37	13	5734	7578	epIZ2380611	mal5 gene product [Saccharomyces cerevisiae] >gpiZ47047ISCCCHR1X_196 Sta1p [Saccharomyces cerevisiae]	33	19
					>pirSIS48478 glucan 1,4-alpha-glucosidase (EC 3.2.1.3) homolog- yeast (Saccharomyces cerevisiae)		
					>gpiM16651YSCS22_1 Yeast (S.cerevisiae) S1 protein gen		

TABLE 3.

Treponema pallidum - Putative coding regions of novel proteins not similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)
1	8	2998	2366
1	9	2391	2651
1	14	6072	7019
1	19	8764	8958
1	20	9778	10317
1	21	9930	10325
1	22	10323	11285
1	23	10832	10473
1	24	11294	10878
1	27	13370	13990
2	2	702	1268
2	3	1027	1950
2	4	2199	1417
2	7	4886	5182
2	8	5095	5343
2	16	13354	13947
4	1	1360	470
4	6	4285	4848
4	10	7150	6740
4	14	9207	8596
5	8	5721	5383
6	1	442	2
6	6	4237	4839
7	1	29	433
8	1	2321	1545
8	2	2491	2105
8	4	2761	3090
8	5	3326	3568
8	6	3566	3751
8	7	4421	3717
8	8	4388	4885
8	9	4785	5204
8	11	6460	5501
8	13	6741	7334
8	14	7189	7557
8	18	9207	8998
8	38	23258	23581
8	39	23766	23395
9	1	1	504
9	2	2312	588
9	3	3895	2333
9	4	3116	3475
10	2	2252	2515
10	3	3473	2904
10	4	4003	3419
10	5	4305	3829

TABLE 3.

Treponema pallidum - Putative coding regions of novel proteins not similar to known proteins

10	6	4439	4107
10	10	6172	6744
10	11	6546	7367
10	12	7175	7546
10	13	7792	7595
10	14	8345	7830
10	15	8248	7940
10	16	8292	9110
10	17	9101	8778
11	1	122	469
11	2	849	1358
11	3	1129	1563
11	4	1624	2919
11	5	2826	3164
11	6	3047	3991
11	7	3940	4557
11	9	5188	4898
11	10	5190	5477
11	13	6392	6111
11	16	8220	7573
11	20	9910	10440
11	21	10917	10402
11	22	10438	11469
11	23	11607	11053
11	24	11950	11576
11	26	13738	14235
11	27	14115	14321
12	1	40	585
12	2	2403	352
12	5	3686	3952
12	9	4715	4957
13	2	680	1303
13	3	1613	747
13	5	2555	1920
13	7	3468	3130
13	9	6157	5279
13	13	8327	8638
15	1	1	759
16	9	6239	5763
16	11	7440	7069
16	12	7901	7194
16	15	9565	10185
16	17	12182	11706
16	18	12584	12300
16	21	12793	13308
16	22	14166	13264
16	23	14450	13992

TABLE 3.

Treponema pallidum - Putative coding regions of novel proteins not similar to known proteins

16	26	16277	15594
16	31	17895	18332
16	32	19003	18242
16	33	19491	18970
16	34	20385	20053
16	40	22895	23161
16	42	24264	24581
16	43	24452	25549
16	45	26020	26673
16	46	26901	27722
16	53	31056	31415
16	55	33441	32989
17	2	1345	962
17	3	2034	1492
17	7	3509	3105
17	8	4602	3892
17	10	5335	4946
17	11	5611	5381
17	12	5645	6409
17	13	6744	7742
17	15	8000	7836
17	16	7839	8450
17	17	8619	8434
18	6	5410	5991
18	8	6421	6020
19	6	4526	4951
19	8	6964	7773
19	9	7544	8032
19	10	8049	8486
19	12	9692	10201
19	15	12883	13215
19	17	13106	13624
19	23	18224	17556
19	24	19216	18200
20	2	1814	2575
20	3	2511	3299
21	13	9403	9813
21	14	9489	9917
21	17	10312	10758
21	19	10602	10946
22	2	1671	2450
22	3	3504	2578
22	4	3372	3869
22	6	5309	4761
22	7	6062	5445
22	8	5592	5975
22	9	6144	7343

TABLE 3.

Treponema pallidum - Putative coding regions of novel proteins not similar to known proteins

22	10	6416	6790
22	11	8006	7293
22	12	7850	8323
22	13	8617	9156
22	18	12328	13269
23	4	1517	1212
24	3	811	1548
25	3	1729	1169
26	2	167	1576
26	3	1512	2843
26	8	4822	4983
26	9	5604	4888
26	10	5855	6100
26	11	6040	6558
26	13	7758	7267
26	14	8349	7903
26	18	9960	10109
26	30	17705	19171
26	31	18835	19236
26	32	19161	19769
26	33	19688	22006
26	34	20812	20591
27	1	669	4
27	2	788	195
28	5	3158	2616
28	6	3367	3071
28	9	4459	4707
28	12	5186	5794
28	15	6636	7073
28	16	7839	7339
28	17	7971	8477
28	18	8389	9942
28	20	9831	10127
28	21	9914	10801
29	3	804	1616
29	4	1519	2061
29	5	1976	2404
29	9	3478	3705
29	12	4800	5204
29	15	6102	5584
29	18	7459	7866
29	21	9443	8706
29	27	12452	13120
30	1	111	1217
30	3	1735	1517
30	16	12128	12544
30	18	13844	14488

TABLE 3.

Treponema pallidum - Putative coding regions of novel proteins not similar to known proteins

34	5	3334	3017
34	6	3095	4246
35	2	856	230
35	3	1968	514
35	4	2807	2073
35	5	3424	2282
35	9	6057	6809
35	11	6867	6592
35	12	7768	7217
35	13	8548	7856
35	17	10143	9238
35	18	9728	10747
35	19	11411	10851
35	35	19000	18431
35	39	21598	21347
36	4	3846	3448
36	5	4694	3759
36	6	4920	5420
36	7	4991	5542
36	8	5540	5842
36	9	6906	5854
36	10	5950	6246
36	11	6816	7121
36	12	7418	6858
36	23	18768	19202
36	25	20011	19688
37	2	1221	751
37	4	1815	1438
37	5	2429	1971
37	9	4315	3716
37	12	5504	6412
37	15	9075	9386
37	17	11158	10835
38	4	1964	1626
38	5	1791	2450
39	2	826	134
39	4	1353	781
39	6	2660	3136
39	7	3336	2962
39	11	6254	5331
40	4	1580	2065
40	5	1807	2034
41	6	2078	1716
41	10	3059	3976
41	11	3865	4854
41	12	5143	4844
41	15	5944	6429

TABLE 3.

Treponema pallidum - Putative coding regions of novel proteins not similar to known proteins

42	2	931	1203
42	4	2149	1514
42	25	18100	17825
43	2	1542	1931
43	4	2842	3246
44	1	351	67
44	3	534	349
44	4	825	532
44	10	5146	5505
44	11	5788	5171
44	15	9156	8251
44	20	12138	12524
45	2	2281	1058
45	8	8147	8656
45	9	8437	8760
45	11	9062	9373
45	15	13884	12790
45	17	14500	15108
47	3	1317	1550
47	5	2545	2105
47	6	2406	3014
47	9	4837	4592
47	10	5173	5700
47	15	8714	9049
47	17	10561	10001
47	19	11772	10807
47	24	13781	13425
47	27	15122	14625
47	28	15917	15480
47	30	16213	16638
48	2	1210	725
48	3	2133	1198
48	5	3031	2483
48	13	9988	9617
48	14	10323	9721
48	15	10538	10149
48	16	11509	10493
48	17	11818	11507
48	20	12934	12656
49	2	1116	823
49	3	866	1435
49	4	1667	2584
49	5	3915	3295
49	6	5194	3806
49	7	5691	5113
49	14	9465	9136
49	15	9467	9760

TABLE 3.

Treponema pallidum - Putative coding regions of novel proteins not similar to known proteins

49	16	10114	9578
50	2	252	611
51	2	3123	3671
51	7	7326	7664
52	1	2445	1756
52	2	2956	2459
52	4	4723	4109
53	1	3	365
53	2	172	603
53	4	1190	924
53	5	1755	1411
53	6	3185	1647
54	16	5423	6022
54	17	5994	6335
54	21	7481	7143
54	22	7243	8028
54	23	7289	7525
54	24	7668	7985
55	2	703	464
55	6	4738	6135
55	7	5937	6731
55	10	8692	9015
56	10	6010	5567
56	11	5699	6061
57	1	3	446
57	10	7281	7565
58	2	1204	551
58	5	2883	2335
59	2	568	867
59	3	2094	1591
59	4	2239	1871
60	2	637	2
60	7	5121	4279
61	2	666	1283
61	3	1263	976
61	4	1281	1778
61	5	1736	1392
61	7	2988	2608
61	10	4819	5673
62	5	3223	2078
62	6	2669	2959
62	7	3389	3099
63	2	553	341
63	3	1189	893
63	4	1916	966
63	5	1729	1992
63	7	3177	2464

TABLE 3.

Treponema pallidum - Putative coding regions of novel proteins not similar to known proteins

63	10	3552	3887
63	11	3926	4531
63	16	6897	5938
63	17	6446	6099
63	18	7228	6995
63	26	14887	14240
63	28	17256	16678
63	29	17509	17033
63	30	17855	17562
63	33	19853	19116
63	35	20299	21408
64	1	236	45
64	2	1029	196
64	6	2960	3196
64	7	3785	3159
64	15	6989	7399
64	16	7273	7719
64	17	7692	8126
64	18	8250	8816
64	19	8749	9192
64	20	9060	9332
64	22	9330	10403
64	23	10052	9453
64	24	10256	10879
64	26	10922	11665
64	27	11590	11970
64	29	12472	12233
65	1	43	378
65	6	2556	3329
65	7	3194	4339
65	8	4105	3431
65	9	4176	4448
66	3	1622	1425
68	3	2093	1686
68	4	2316	2119
68	5	2399	3049
68	6	3013	2450
68	8	3390	2881
69	1	3	743
69	2	629	886
69	3	1397	906
69	4	2091	1342
70	4	2237	2662
70	8	6703	6996
70	9	7148	7954
70	10	7671	8378
70	12	8925	9299

TABLE 3.

Treponema pallidum - Putative coding regions of novel proteins not similar to known proteins

70	13	9103	9369
70	15	10187	11269
70	16	11447	11022
70	17	12491	11793
70	19	13758	13964
70	24	16029	15769
70	27	17085	17675
70	29	18387	17896
70	31	19514	18957
70	33	20556	20128
71	3	1350	1105
72	5	3863	4087
74	1	129	611
74	2	676	380
74	3	440	916
74	4	820	1728
74	6	2031	1285
74	7	1820	1545
74	9	3676	4929
74	10	4994	6424
74	12	7476	7748
74	13	7801	7634
74	14	7640	8152
75	1	748	2
75	2	1221	655
75	3	2151	1882
75	4	3249	2212
75	5	2233	2559
75	7	2932	2375
75	8	2683	3639
75	9	3507	4034
75	10	3777	3559
75	12	4514	5227
75	13	5377	4922
75	14	5861	5556
77	3	1365	2441
77	4	2317	2832
77	5	2514	2888
77	6	2816	3664
77	7	3517	3912
77	9	4475	4855
77	10	4788	5060
79	6	2455	2727
79	7	3038	2772
79	13	7889	7512
79	15	8631	8059
79	16	8802	9056

TABLE 3.

Treponema pallidum - Putative coding regions of novel proteins not similar to known proteins

79	18	10273	9752
79	19	10092	10436
80	3	1228	551
81	2	1024	611
81	3	1970	888
81	4	1605	1997
81	10	7615	8082
81	16	14381	14085
81	18	15195	14815
82	1	3	260
83	1	1	897
83	2	351	866
83	5	1759	2667
83	6	2565	2984
83	7	2932	3282
83	9	3362	3784
83	10	4123	3494
83	12	4452	4105
83	17	7929	9203
84	5	1484	1813
84	6	2019	1810
84	7	2829	2377
84	8	3531	4421
85	1	1	510
85	3	762	1235
86	4	2070	1549
86	5	1600	2073
87	3	2608	2021
87	4	2615	2274
88	2	894	2081
88	3	2625	1975
88	7	3403	3750
89	2	992	1432
92	2	660	1043
93	1	17	580
93	2	626	1348
94	3	1869	1480
94	4	1706	2074
95	1	2	184
95	2	37	327
95	4	1018	1494
95	5	1685	1284
95	6	1881	1624
96	3	586	191
98	4	4094	3825
98	6	4968	4618
100	2	1356	4

TABLE 3.

Treponema pallidum - Putative coding regions of novel proteins not similar to known proteins

100	3	653	183
100	5	1095	1532
100	6	1283	1585
100	7	1450	1836
100	8	2010	1459
100	10	2149	2619
100	13	3049	3573
100	16	3631	4023
100	23	7950	7630
100	24	7690	8220
100	25	8201	7890
100	26	8592	8990
100	27	10023	8770
100	35	14886	15329
100	36	15272	15769
100	38	16962	17621
100	40	18590	18327
101	3	802	515
101	18	8545	9162
101	22	13204	12545
101	24	13821	14135
101	27	16471	16788
101	28	17178	16543
101	29	18178	17063
101	30	18212	18847
101	32	19052	19798
102	2	1415	921
102	4	2083	2730
102	8	5114	4851
102	9	5420	5112
103	3	1699	2436
103	6	4533	5243
104	2	473	1069
104	3	1085	1711
104	4	1411	1695
104	6	1896	2405
104	8	3705	3145
104	9	3729	3983
104	10	4782	4192
104	14	5889	5557
104	16	6123	6404
105	1	472	68
105	2	890	375
105	3	849	1592
107	1	1575	898
107	2	2213	1503
107	3	2994	2173

TABLE 3.

Treponema pallidum - Putative coding regions of novel proteins not similar to known proteins

107	4	4979	2892
109	2	1429	497
110	2	698	1399
110	3	1396	1133
111	3	1854	1486
112	3	2233	1868
112	4	3078	1915
112	13	9340	9005
114	2	576	190
115	3	1551	1090
115	5	2499	1933
116	1	346	5
116	2	439	1143
116	3	801	478
116	4	1254	919
116	6	1646	1188
116	8	1968	2459
116	11	3583	4077
116	12	4235	4852
116	14	5140	6093
116	15	6922	6170
116	16	6483	6722
117	8	6315	5743
117	10	6984	6493
121	2	604	1314
121	9	3632	4156
121	10	4524	5723
121	11	6081	5785
121	12	6015	6359
121	13	6598	7308
121	15	7894	7637
121	17	8347	8027
121	20	9966	10946
121	21	10457	10110
121	28	16019	15696
121	29	17161	16733
121	30	16754	17380
121	31	18342	17401
122	1	325	729
122	2	603	1808
122	4	2300	2785
126	1	191	3
126	3	906	724
126	4	1484	2170
134	1	486	283
135	2	880	554
135	5	2797	3129

TABLE 3.

Treponema pallidum - Putative coding regions of novel proteins not similar to known proteins

135	8	3381	3695
137	1	2	970
137	2	180	22
137	4	1037	1321
137	5	1319	3001
137	6	1966	1601
137	7	2707	3690
137	8	3501	3169
137	9	3414	4613
137	10	4726	4535
137	12	5477	5160
139	3	819	1229
140	3	1465	671
141	2	1273	848
141	7	4071	4331
141	9	4761	5282
142	7	10511	9729
142	8	10514	11344
142	9	11325	12233
142	10	13391	12678
143	6	5906	5304
143	7	6546	6130
144	3	640	1065
144	5	1627	1121
144	6	2102	1593
144	7	2704	1979
144	8	2898	2344
145	1	3	983
145	2	728	51
145	3	763	341
145	4	890	2065
145	5	1908	1570
145	6	1927	2478
145	7	2499	2011
145	8	2304	2669
145	9	2629	2967
145	10	3141	2638
145	11	2877	3290
146	1	3	1316
148	1	1	261
148	2	1145	282
149	1	16	489
150	4	3743	3285
150	10	8189	8986
150	12	11486	10512
150	13	10524	10994
150	15	12079	12366

TABLE 3.

Treponema pallidum - Putative coding regions of novel proteins not similar to known proteins

152	2	2411	1977
156	3	561	139
156	4	883	635
157	1	1	762
157	2	662	1231
157	4	2981	1695
158	2	456	43
161	1	466	5
161	3	558	857
163	1	607	140
163	2	716	468
164	1	1283	3
165	2	369	88
169	2	343	564
172	2	609	295
172	5	1978	1466
172	6	2076	1792
172	7	2825	2019
172	8	2424	2864
178	1	98	556
182	1	248	3
186	3	734	1267
186	4	1113	1379
188	1	3	686
188	2	310	843
189	1	688	2
192	1	245	3
192	2	18	413
193	1	85	507
199	1	146	376
203	1	321	611
206	1	2	568
209	1	1	543
210	1	229	2
212	1	42	584
212	2	383	808
224	1	38	286
224	2	579	325
276	1	201	587
328	1	360	4
376	1	567	139
389	1	485	3
423	1	545	270
478	1	277	11
480	1	27	305
482	1	327	79
484	1	310	8

TABLE 3.

Treponema pallidum - Putative coding regions of novel proteins not similar to known proteins

522	1	473	75
524	1	248	3
551	1	8	547
558	1	3	455
559	1	537	55
565	1	82	420
566	2	360	929
566	3	769	1104
579	2	379	2
605	1	334	53
608	1	186	4
620	1	444	115
625	1	281	3
626	1	253	41
626	3	847	578
628	1	555	79
628	2	626	306
633	1	195	4
634	2	35	583
636	1	3	308
643	1	1	402
644	1	1	339
644	2	525	4
645	3	747	427
646	1	79	453
648	1	426	4
649	1	264	536
659	1	90	359
668	1	103	342
668	3	288	536
669	1	251	39
678	1	382	95
679	1	513	130
682	1	108	434
684	1	438	133
687	1	2	262
691	1	337	14
702	1	549	121
703	1	2	307
719	1	531	358
742	1	408	220

(1) GENERAL INFORMATION:

(i) APPLICANT: Human Genome Sciences Inc., et. al.

(ii) TITLE OF INVENTION: Treponema pallidum Polynucleotides and Sequences

(iii) NUMBER OF SEQUENCES: 744

(iv) CORRESPONDENCE ADDRESS:

(A) ADDRESSEE: Human Genome Sciences, Inc.

(B) STREET: 9410 Key West Avenue

(C) CITY: Rockville

(D) STATE: Maryland

(E) COUNTRY: USA

(F) ZIP: 20850

(v) COMPUTER READABLE FORM:

(A) MEDIUM TYPE: Diskette, 3.50 inch, 1.4Mb storage

(B) COMPUTER: HP Vectra 486/33

(C) OPERATING SYSTEM: MSDOS version 6.2

(D) SOFTWARE: ASCII Text

(vi) CURRENT APPLICATION DATA:

(A) APPLICATION NUMBER: Unassigned

(B) FILING DATE: June 23, 1998

(C) CLASSIFICATION:

(vii) PRIOR APPLICATION DATA:

(A) APPLICATION NUMBER: 60/050,667

(B) FILING DATE: June 24, 1997

(viii) ATTORNEY/AGENT INFORMATION:

(A) NAME: Brookes, A. Anders

(B) REGISTRATION NUMBER: 36,373

(C) REFERENCE/DOCKET NUMBER: PB387PCT

(vi) TELECOMMUNICATION INFORMATION:

(A) TELEPHONE: (301) 309-8504

(B) TELEFAX: (301) 309-8512

(2) INFORMATION FOR SEQ ID NO: 1:

(i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 14063 base pairs
- (B) TYPE: nucleic acid
- (C) STRANDEDNESS: double
- (D) TOPOLOGY: linear

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 1:

AAGGTTGTTT	GTTGATAATC	TCTGCCATAT	TACTGTTCCC	TTCTTTTCGT	TCATCGGGTA	60
AGAGCCGTCA	GCGGTGAGCG	CGCCACcTCC	TTCACTAATC	ACCACGTCGT	GAACAGGC	120
TGCcGTAcAG	CCGACCACCA	ACGGCTCTTC	TACCCACTCA	ACTCTACATT	CCACGCTAGC	180
GAGCCAACGC	AGCAACGATC	GTATCGATAT	TGTGTGTTAC	CATCCCTACG	TAGGTACCC	240
CGCTCGTACC	CGCATCCCCC	ATCGCATCAG	AAAACAAC	TGCnTACGTG	300	
CCCTCTTGCC	TGCACCGCAT	CCCTTAACGC	TTCAACGTTT	TTGTGCGGAA	TAGAACTCTC	360
AATAAAGATA	GCAGGGAGTT	nTACnTGCGC	AATAAACGCT	GCCAGTTCCCT	GCATATCATG	420
CGCACTGGCT	TCCGAAGCGG	TGCTCACCCC	TTGCAACCCC	TTCACCTCAA	AACCATACGC	480
ACGGCTAAAA	TAGCCGAACG	CATCATGAGC	GGTCACCAAC	ACACGCcTTT	CAGCAGGCAG	540
CGACTGCGCC	TTGGCCGAA	CGTACGCGTC	AAGCTTATCC	AACTGCTGCT	GGTACGCGCTG	600
ATAACGTTGA	GTAAATTGCG	GAGTTTTCC	CGGCAACAGC	TTGCACAAGC	TTTCGTACAC	660
TGCCTTCACC	GAATAAGACC	ACAGCTTTAC	ATCAAACAC	ACATGCGGAT	CGAACTCTGC	720
TTCCTCAAGA	GAAAGACGCT	GAGACACCGG	AATAGTCTCA	GAAACTGCAA	CTACCAAGCG	780

179

GCTCCCGCGC AGTTGGAAA ACACCTCGCC CATCTGGTT TCCAGGTGCA ACCCGTTGTA	840
CAGGATGAGA TCCGCATTCC CGAGCCATTG CACATCCCC GCAGTAGCCG TGTACAGGTG	900
CGGGTCAACA CCAGGACCCA TCAACCCCTT TAGATGCACA TCACCTTGAG CGATGTTTT	960
GACAGCATCC GCTATCATGC CAATGGTGGT GACAACCAGG GGTTTCCCGT CCGCTGCGGC	1020
ATCCTTGCTA CCGAATGCGT GCGTAAAACC GGTCAGCATG CCAAGCGCGA GCACGCAGGC	1080
ACATATTCTT TCACGTATCA AGTGACACTC CTTGGGTGAA TTTgATGCAT CAAAGTAGCG	1140
GATCACGCCAG GAAACGTCAA TTTTCAGTA CCCCTTCAGG AAAAGAAAAC GGCACCTCTGG	1200
CGCGGACCTA CTCGAGCTAC ATGATAAAGA AGCATTGAT CTTCCCTCGCA TAAGCAGCTA	1260
CAGTTGCGCC CTCTTATGGA TCACACCGCG TCACTGAGTC CTGTGCGCCC TGAAGCACAA	1320
CCTACAGACG ATCGCGAGCG TGCGCTCAGA CCGGCCCTCC TGAAAGACTT TCTAGGTAG	1380
GAGAAAACAA AACGCAACTT ACGTCTTTTC ATTCAAGGCAG CGCGCGATCG CAACGAAAGC	1440
TTAGATCACC TGTTCCTCAT CGGCCCCCGG GGGCTcGGCA AAACGACGCT CGCGCATATC	1500
ACTGCATGCG AGCTGGCGT TGAGTGCAAG GTTACAGGCG CACCGCGCT TGATAAACCA	1560
AAAGATTAG CGGGTATCCT CACTGCGCTG AGTGAGCGAA GGChTTCTTC GTGGATGAAA	1620
TCCACCGCCT CAAACCAGCC ATAGAAGAGA TGCTGTACAT TGCCATGGAG GACTACGAAC	1680
TGGATTGGGT TATCGGTCAAG GGACCGTCCG CGCGCACGGT GCGCATCCCA CTCCCCCGT	1740
TTACCCCTCAT TGGTCAACC ACTCGCGCGG GTATGGTTTC AAGCCCGCTG ATTAGCCGCT	1800
TTGGAATCGT AGAGCGCTTC GAGTTCTATA CCCCTGAGGA GCTTGCTGCC ATTGTGCAAC	1860
GCTCAGCGCG GCTTCTAGAT ATCACGCTCG ACGCACCGCG AGTThAGCCC TTGCGCGGTG	1920
TTCGCGAGGA ACACCCCGGG TGGCCAACCG GCTTTTGCGC CGTATAACGCG ATTTTGCCCA	1980
AGTTGCGGGG TCTGCACACA TCAGCGAGAC GATAGTACGC GCAGGcTTGC CCACCTAAAG	2040
ATCGACGAAT TAGGGCTAGA ACTGCACGAC ATACAGCTGC TGCGCTCAT GaTTGAGCAC	2100
TTCGGCGGAG GGCCAGTGGG CGCAGAAACG CTGGCGATCT CCCTCGGGGA ATCACCGGAA	2160
ACACTTGAGG ATTACTACGA GCCCTACCTT ATCCAAATTG GGCTCATGCA GCGCACCCCC	2220
CGCGGGCGCA TGGCCACCGC GCGTGCCTAT GCGCACCTAG GTCTCCCTGT CCCCAGGGCA	2280
CGCACGCTCA CCCCCGCACTC CCCAGAACAA GGAACGCTTC TTTAGCAAAG ATGCGGACAC	2340
CTTGTCAAG AGCTGGGAGGG GTGCGTTTCG TGAAGATGTC TGCCTTTTT GCACCAACCT	2400
ATATACTACG CGCGGGAGGG GTGCGTTTCG TGAAGATGTC TGCCTTTTT GCACCAACCT	2460
GCAGGCTTGCA CCTGCTGATG CAACCATCGC AAGCCACCGAG CTGCTCATGC GCGCAGGGTA	2520

180

CGTCAGAAAA ATCGCCAACG GCCTGTTGCGTACCTTCCC CTGGGCCtGCGCTGACA	2580
CAAAATTGAA GCGATTATTC GGGAAAGAACT CGAGGCTATC GGGTGTGGAAGTGCACCGC	2640
GCCTGTCGTG ACTCCTGCAG AGTTGTGAA GGAATCTGGCCGCTGGTACCGCATGGGCC	2700
AGAGCTTTG CGCGCCAAAA ATCGGCTCGA TCACGAGCTC CTTTCAGTC CGACTGCAGA	2760
AGAACCTTC ACCGCTTTGG TCGCGGGCGA CTGTACTTCC TACAAACATT TTCCCCCTCAG	2820
TCTCTACCAA ATCAACGCAA AATATCGCGA TGAAAATCCGT CGCGCTTACG GACTGATGCG	2880
CGCGCGCGAG TTCACCATGG CCGACGCCTA TTCTTCCAC ACAGACTGCG CATGCCTTGCG	2940
GCGCACGTAC GAAAAGTTTG CCCACCGCGTA TCGCGCCATT TTCCGTCGCA TCGGCCTATC	3000
AGTCATTGCA GTACATGCAC AccTCGGTGC GATGGGGGG CAGGAATCCG AGGAATTGAT	3060
GGTAGAGTCC GCGGTGGGCG ACAACACGCT CCTGTTGTGT CCCCCACTGCA CCTACGCTGg	3120
CAAATTGCGA AAAGGCCGTC GGACAGCGCC CCCTCCCAGA CACCGCATGAC ACTCATCTAA	3180
AAGACGAACA CGAAGGgTCA GATCTCAAGA CGCCTGCAGC AATGCGCGAG GTGCACACCC	3240
CGCACGTGAA AACTATTGAG GAACTTGAAAC ACTTCTTGCA CGTACCTGCA CATCGCTGCA	3300
TCAAGACGCT TATTTACCGC ATTGACACGG TGCCCCAGGC GGCTGGGCAT TTTGTGGCAG	3360
TGTGCATCCG CGGGCACCTA GAACTCAACG AGTCAAAGCT CGAACGGCTC CTGGCGGTGCG	3420
CATCTGTAGT ACTGGCAACT GAACAAGAGG TGTATGCACT CAGCGGCACC CCCGTAGGAT	3480
TCATTGGTCC GGTAGGACTT GCACAGCGTG CTGCAGCTGC GTATGCGCT CGCACCCtGC	3540
GTTCTTCCCC TCCGCTGCTG AGCCTGCATC CGTCACTTCT GACATTCCAT TTTTTTCCCT	3600
CGTTGCAGAT CAGTCCGTGA TGGCTATGCA CAACGCTATC ACCGGTGGGT TGAAAAGTTGA	3660
CACGCATCTT GTGCAGGTAG AACCGGGTCC AGACTTTGTT CCTGACGCAg TTGCAGATCT	3720
CATGCTCGTG CGCGCCGGCG ACCGGTGCAT ACACGTGGA GCGCCCCTAT ACGAAAAAAA	3780
GGGTAACGAA CTAGGTCACC TCTTTAAATT AGGGGACAAA TACACGCGCA gATGcACCT	3840
TACCTTTACT GATGAGCAGG GTGTACGACA GTPCCCCCTG ATGGGCTGCT ATGGCATTGG	3900
CCTTGATCGC ACGCTTGCT CTGTGGTGGAA AACCCACCAT GACACGGGG GTATCAGCTG	3960
GCCGCTTGCG ATCAGCCCCCT ATGCAGTTGT GCTCATACCC ATCCCTCACA CGCAGGGCCC	4020
CTATGCAGCA GCAGAGGCAC TGTACGTGCA GCTGCGGACA CGGGGAGTTG AGGTACTGTT	4080
TGATGATCGT GCAGAGCGAC CCGGAGTAAA GTTCGAGAC GCTGATTAA TCGGTATTCC	4140
CTTCGTGTGG TACTGAGTGC GAAAAAnCTAC CGCGCGTTGA ATGcaCAACA CGGTGTGGTG	4200
CGCACACGTA TTTTTTACG CAAGAAGAGG CGTCCGAGCA CATTGCACGC CTGCTCGAAC	4260

181

AACTCGCTTC CCCGGAAAAGT TCGTAAGAAC GGGATGCCG GAGCGGGATC CAGCGCATGC	4320
AGTGCTGAGA CCTGCGCATA ATAGCACAGT GTACGGCACC CGTGGTTTAG AAAAAAAATGA	4380
CGAAGGAGAA AAGGAAAAAC GGTGTACATA AAGGTAGCGC TCGTGTGTCT TTTCAGCATG	4440
GGAGGCCGGT GTCTTTGGC CACAGAACCG GCGCCAGTCT CTGGAGATT ACGTATTGTAT	4500
CGCGACTATT CGTGGAAATC GCCCACATGG GTTGGCTTT TGTGCTACGA CGCACACACG	4560
TACGGTGCAG TGCTGTGTAC TCCGGCAGAA AGCCGCAGGA TCACAATTCT CTTCACGGGT	4620
ACTGAAAAGC ACGGCCGCTT TGAGCTGACC GGACAACGCA TCACCTCACC GGTGCGCACA	4680
GAGGATCTGA CTGGCATAAA TTATCTCATG GATCTTTTC CTCAACTACA GCGCTGGAAG	4740
CATTTTCCCC GGGATACACA CACCCCTGTT GCGCGGCATA CCGATCGGAG TAAAAAGAGC	4800
ACACAATTCT CAGGGGCAGT CGAACTGCAG TTCGCTCTT TTGTCCCCCT CTTCCACCTA	4860
GAAATACTCC GTGATAAGCA GCAGCGCGTC ATGCTCCAGC TAAGCGAGAT AGGGAAAGATC	4920
GACCACACCA GTGACGCAGC CTTCTTTCAA TTCACCCCCA TGCCCCCGTC CACGCCCACT	4980
GATGCACCGC CAGCAACGCT TAATCAGACC CTGACACGCA CGGAGTATGT CATCGATGAC	5040
GTGTGCATTG CACTTGATCC GCAGTGGAAA AGAATTGCAG AAAATTCTTT TCTTTCAGAC	5100
TTTGCCTTTC TCACCGTACA CCAGGTGCCT GCACCGCGCG CGCACGACTA TTCTGCGCTC	5160
CGTGCATTGC TGCAACTCTT TCTGTATTCA GGCCCTCAGG GAAAAAACAT TCTTGAACAA	5220
CTCCATATCA ATGACACTCA CGCGCGCTT ACGCTTTCCCT ATGCAGTGTGTT TGACCTTCCG	5280
TCAAAAACAG TTAAAAAGAC ATGGAAGATA TTCATCCGCC ACTCTGATAC GCACTACTCT	5340
ATACTTAGTC TCACGGCGGA CCAgCGCACA GGGCAGSGTT ACGCGCGCTA CTTTGACACG	5400
CTCATTGAAA CTATCCGTAC AAAAAACTAA AAAATGCTGA ATTGGAGCAT ACCCGTGATT	5460
AGACACATAT TATTGACAT AGACAACACG CTGACTCCCT GTACAAATCC CATTGAAATG	5520
GCTATCACGC AGCGCATACA CACATTGTT GCACATTTC TCCACGTATC TTGTGAGGAG	5580
GCGCGCGCGT TACGCCAGCG CACAAAGCAC CTCTATGCTA CCACCTTGA GTGGTTAAAG	5640
GCAGAGCACA ATCTCATTCA CGATGAACAC TACTTCGTG CCGTATATCC TCCCACCGAA	5700
ATACAGGAGT TGCAGTACGA TCCGATGCTC CGCCCTTTTT TACAGTCACT GCACATGCCA	5760
CTGACGGCAT TAACTAACGC ACCGCGCGTG CACGCACAAAC GCGTATTGGA TTTTTTTCAT	5820
CTGTCAGACC TTTTTTTAGA TGTCTTTGAC ATCACGTATC ATGCAGGCAA GGGAAAACCA	5880
CACCAACAGCT GCTTTGTACG TACGCTTGAA GCGGTACACA AAACATGTGCA GGAAACGCTT	5940
TTTGTGATG ACTGTCTCAT GCACGTGCGT GCCTTTATTG CGCTTGGCGG ACATGCCGTG	6000

182

CTGGTTGACG AACGTGACTG TCATGCAGAA CTGCCCTT CTGCACGCAT GACACGCGTA 6060
 AAAACAATTT ATGAATTGCC CGCACACCTT GCACGCCTCG CCCAAGGAGA CAATCAGTGA 6120
 GTATACATTC GTTGCAGCAG ACTTTTAGCG ACATCGTCCC GCTCCTGGAG CAGTATAACGC 6180
 GCGCAGACCG CTTCATGCGG GAGGATAATT TGTTACACGA GAGAAACGAA CCTATCCGGC 6240
 GTATCGTTGA GTCCCTCGTC GCCCGCATAT TACTCCCCGG CTCCACAATG CGCGGAAATG 6300
 AGCAAATCGC ATCCTTTTA CATAAAACCA ATGAAGGGAA ACGGGGACTC ATTCTTGCGG 6360
 AACACTACAG CAATTTGAC TTACCTGTC TGCTCTACCT TATGGAACAA GGAAGTAGTG 6420
 CCGGGCGCAT GCTTCAGAA AAAATCGTAT CTATTGCCGG TATTAAACTT CGTGAAGAAA 6480
 ATCGCATCCT GGCAATGCTC ACCGAAGgAT ATGATCACCT GGTGATATAT CCCAGTAGGA 6540
 GTTTGGCCAC CATCACTGAT GCGCACTGTC TTGCAAGAGA GACAAAGCGC AGCnGAGCAC 6600
 TGAATCGTGC AGCTATGAAG TATTTAGAGG AACTGCGCAA CGCGGGAAAG GTGATTCTCG 6660
 TGTTTCCCTGC AGGGACACGC TACCGACCCG GGAGACCGGA AACAAAGCGA GGGGTGCGCG 6720
 AAGTATACTC CTACATAAAA CACGCCGAGG TACTGCTCCT TATTTCAATC AATGGGAATT 6780
 GTTTGCGCGT TGCAGAACGT TCAACTGATA TGACGGAAGA CGCGGTGCAT CCGGATGTCG 6840
 TGCTTCTTGA AGCCGCACT GTAGACGAcT GCGCCCTTT TCGAGAAAAA GCGCTGGACT 6900
 GGCACCGCAC ACACAACGTG GCGGCACCGT CAGAGGATAA AAAACAAATC GTAGTCGACT 6960
 ATGTCATGCA CCTTTTGAA GAAATGCACG AGCACAAATGA ACGAGAAAGG CTATCGTGA 7020
 TTTTTCGCTG GAATTCCCCG TAAGATCCTA TGAGCTAGAC GGATACGGAC ACGTGAACAA 7080
 TGCGGTATAT CTCCAATATT TTGAATATGC GCGCGCCGCT TTTTGTCTCC ACATAGGGTT 7140
 CGACCTCAAA CAGTTGCACG AAGCAGGTTA CGCTTTCTAC GTAACCCAGG CGCACATTCA 7200
 CTAcCGCACT GCAGTGCATC TATTCGATAC GTTGCGCGCC CGGGTAAAAC CATTAAAGCT 7260
 CGGAAAAGCT TCGGGCGTCT TTTCACAGAC GCTGGAGAAC CAGCATCACG TGCTATGCGC 7320
 GGATGCGGAA ATTACCTGGG TGTGCGTTTC GCGCACAAAGC GGCAAACCAA CTAAGATTCC 7380
 CCCCCGAGTAT CTGGTACCTG CGCTGTATCC GAACTACTAG TCCTCCCTTC TTTCCCCTTT 7440
 ACTCTCCCAA GGACATCACA CTACGGAAGG GTACGCATAC GCAGTAGGGA GGTAGGGTTT 7500
 ATCGCGGAGC CATTCTTATA GATTGTAAA TGCAGGTGTG GTCCCGTGCT GCGTCCTGTT 7560
 TTTCCCAATA ATCCGATTTT tGTCGCGCTG GTGACCGCGC TACCTGCTGA AACCAACACC 7620
 GTCTGCAGAT GCCCCATACAG GGTCTGATAAC CCCGCGTGGT GCCCCACAAT CAGGTAATT 7680
 CCATACACTG CACTGTATCC AACCGTGCGT ACAATCCCTC CGAGCGCCGA ATATACTGGG 7740

183

GTACCCGCC GACTCACCAT ATCCAAACCA TTGTGAAAAC TTCTGGCACC GGTAAACGGA	7800
TCAcTACGCC ATCCATACCG CGAAGAAAACA TAGTACCGAC TGCGAAGAGG AGCACGAAAC	7860
AAGTCACCACAT TAATTTCCTG CAACGCGCGT GCGCTTAAAT GTGCACCGGG CAAAAACAGT	7920
ACGCGTGCAG GCTGCAATGG CTGCAcTGCG TCAAACGACG TATTTCCCT CCACTGTTTC	7980
GCAGAAGAAA ACGGAAAAGG CACGCAGgAC TCCC GTGCAG CTGAATTATA GAACGGAGAA	8040
ACCAGCGTAC GCAC TGAGG AGGTGACTCC TTTGAAGAAG ACGGCGTGTT AAGCAGCACC	8100
AATCGTTCTA AGGAGATCTG ATGCGCCGCC GCTATAGACG AAAACGTATC GCCGTTTTT	8160
ACGGTATATA AAATGCCGTC CACTGAGGGG ATTTTAGTA GCTGTCCAAC TTGGAGCGCC	8220
CGTtGyTGCG CAATTATTTC AACTAATGA TTGCATCCTG ACTGATGTCA TAgcgCtGCG	8280
CAATCCTTCC TACCACATCA CCTTCACGCA tTCGTACACT GTGTAGTACA GTGCAGGcTC	8340
CGCATCTTCC TGCACGATAC GTGCACGGAG CAAGGAAGAC ACGTACCCCG ACGCCTGACG	8400
TGGTTCCCTGC TCAGTGAGCG TGAGGGCAGG TGTCAATGGT TCCACCTGAG CACCAAAGTA	8460
CGCAAGGGCA AGAGCAAGGA GCAACAGTGT TACGAACAGT AACAGTnGCC tACGGGGACA	8520
GGTCTACACG GTTCTCGCAC AGTCTGTTTG GAACTTCGAC AGTACACGCT CACACCGGCT	8580
ATCCTTCAGG TGTACACACT GCCGTATCtG CGGGCAGGTT GCGTCTGTAC CTAACGCACC	8640
GTCTAGAGCG TCCACGCACG CAcGCGCGCG CGCGAGGGAG TCCGGCGGAA AAGAAGTTAA	8700
CACCCGCATG AATGCAGGGC TCGGTGTGT CCACACCTGT GCAGGCAGCG CCTGAAGACG	8760
CGATGAAGGA AAACGCACAC ACCAAGCAAG CAAAAAAAAG nGCGTGTAAAC GCGCATTCCC	8820
GTGcTGACTC GCGCCACCGT ACGTGCATC TCCTACCAAGG GGGATCCCT GTCCAGCGCA	8880
ATAACGGCGA ATCTGATGCT TTTTCCCCGT AACCGGCACA ATCACGCGGA GCACCAGCGC	8940
GCTATCACAG CTATGTAACA CTGTTTGCAC ATGCGTTACC TCTCCTGGAC GCACCAGCGT	9000
GCGCGCCCGA GCAGCGGTGC gCGCAGGGGC GGCGGTGATC GCAAGATAAA ACTTGCGCAA	9060
TGTATGCTGC TGCAACGCGG CAGAAAACCA CTGGGCACCG CGTAACGAGC GCGAAAAAGC	9120
AATCAGTCCC TCTGTCCCTC GGTCCAAGCC GTGCAACGGT CCAGGGCGGA ATGACAAAGC	9180
AGGGGAAACG TGCGCACGCC CTTGTCCCCCT CACCCAGGCA TCCAGGctGC GCGGACCGTG	9240
CACAcAcAAC tGCGGGTTTA TGAAAAAAAAA GCAAATCTTG TGTTTAAAT ACCACCGAAG	9300
cAACACCGCGC ATTcCGTGTT CCAGGCATCT TCGAAAGACG ACTGGATGCT GCACGCGCCC	9360
TACACAGGGA TTCAGGTAAA GAAAGCACAT CCCCCACCTG CACCCGCTTT GCAGGCTGCA	9420
CCGGACGACC ATTGAGCCGG ATAGCGGTGC GGCGCACGCG GGCATACACC CCAACACGCG	9480

184

GACAGGCAGG CAACAATATT CGCAAAACAC GATCTACTCG TCTACCTGCA TCGTTTTGG	9540
TGCAGCGAAA ACACCTAACAC GCCGCTCCAC CATGGGGTCT CACGGTGGAA ACAGGAGGGA	9600
CGACATCCAT ATGCACAGTG TGGGGAACGT TAGACGAGAC CCACCTTTTC ACGCGACGAA	9660
CACCTACTTT CATAACACGGT GCGTCCCCGT GCCGGATACC AGTTGCTCCT CCCAAACGTC	9720
CTCCCCGTCT TTCCCAGTAC GACACCACAG CCCATGGCGG GTACAGCCGC CCCAGTATAG	9780
CGCACACAAAC GCTCCCTGGA CAAAGGTTTA GAGAGTATAG GAGACTGCC CGCGATGGAC	9840
GGTGGCTATT TTCTTGGCCA GCTGCATGCG GTGTTCAAGTG GTGAAGTCTT CCTCTCTGCC	9900
ACCTGTAGTT GGCTTGCAAG TCAGGTGATT AAAGTGGCTA TCGCATGCCG AAGtCGGCTA	9960
TACGGTCGGT GCACGGCTTT TTTGATTTTG CTGTTGGCG CACCGGCGGC ATGCCCTTCGA	10020
GTCACACTGC TCTTGTGTCG GCGCTCACGC TCTCTTTGTC GCTCAAGTGC GGGTTGCATT	10080
CGGATCTGTT CATCTTTTCC TTTTCTCTG CCATCATTGT CGTGGCGCAG GCGCTCGGTG	10140
TGCGCCGTTA AAGCGGCCTG CAGGCCGAGG CGCTCAATAG CCTCGGTGCG CGTGTTCGG	10200
AGAAAACTTGA TTTTTCTTTAG ACCAGTGC GAGAGATTCA TGGACATAAA CCGCTGGAAG	10260
TTGTCGTTGG CGTGGCAGTG GGCAATCGTCA CGAGCGCTTT GTTCTACAGC TCCATGAGCC	10320
CTTGAGTCTC CGGTGGACGT GCATGCAATG CGGcGGACCC CTCCACACAG AGGAAGAGGC	10380
GGTGCCTGTC GCGTTCTCTC GTGTGTCCCT GCGCGTCGG GAGGCCGAGA CCTTTTCTGT	10440
ACCGTACAGA GGGCACACCA ATGATAGAGC GCCTACGGAG CAGTCGCGGG AAACTCACCC	10500
TCACCCACCA GATTTTCCCC CTCAGCTTTG GGGGAATGC TTTTTGCGCT GCGCGCGCG	10560
TCGTTCCGTT CTCCGTTGAT GCTGGAGAGC CAGCCGCCGT CGCTGTGGTA AAGGTTGGGG	10620
ATACGGTCAG AGAAGGTCAAG CTGATCGCAC GCGCCGCGCA CGCCGGTGCT GcTCACGCAC	10680
ATGCCTCCGT CCCCAGGTGTC GTCACCCGCT TGGTAAGTGC TAATTTCTC GCCGGTAGTG	10740
CCCTGCGCGC TGTCGAGATT CGTACACGCG GTTCCTCGA ACATCTTGGC AAGGTCAAAC	10800
CAAATCGCCC GTGGCAGCAC AGCACCGCTT CAGAATTGct GCGCCTAGTT ACAGATGCAG	10860
GAGTAGTGGC CACACGCCA CATCCGCACG CCCAGATCAC GAGCACCGCA ACGGGCACGC	10920
ACGGGGTGC ACAGCACACG TACCGGAAAG ACTACGGACA GAAGAGAAGG GCTGAAGCGC	10980
ACACGCTGCG TCTCATGCGC GCGCGTGCGG AAAGCGGCAA TGCGCTGCC ACACCGCTCC	11040
ACCTGCACGT GCGTAAGGGT GTACGGAAAC TTACGCTCTA CCTTTGTGAC GACGACGCTA	11100
CCTGCCCTTT GAGTTGTTTC CTTGCGCAGG AGTTTCCAGA ACCTGTTGCT ACCGGTACCG	11160
CCATTATTGC ACGGATACTG GACGCTACGT ATACCCGCGT GTCTCCACAC GCTGCCAAAA	11220

185

CGCTCCCCCG GTCTTGCAAG GATGCGCGCT GTCTTCCAT TCAACGAGAT GCACGACGCG 11280
 TATAGACGAC ATTATCCTTT TAGCAATCTA TGTGCCAC GCTATCGTGC AGGTTGCACA 11340
 ATCGATGCAC TCACTGCAGT GCACGTGTAT GAGGCAGTGG TACTCAGTCA GCCGCAAATC 11400
 AGTT CCTACA TTGCTCTGAC AGGCAGTGG A TAAAATCAC CGCAGGTACT CCGCGCGCGT 11460
 ATCGGCACCC CCCTTGGCGC GCTCATCGAG GAGTGTGGAG GGTTTCGCAC ACGCCCCGGG 11520
 CATCTCATCA TCAATGGACT GCTCAAGGGT AGTGTGTTAG AGTCGTTGGA CCTGCCTTTC 11580
 TCAAAGGGGA TCAAATCGCT CCACGTCACC GGTAAGCGC TTTCAAGCTC TGCGTCCGT 11640
 ACCTCCTGTC AAAACTGTGG TGATTGCGCG CGCATTGCC CAGTATATCT TGACCCAATA 11700
 AAAATTGCGC GTGCCGCACA CCGTAATCAG TTTACTGAAG AAGTGTCCA ATCCcTGcGG 11760
 ATTTGCCACC AATGCGGTCT GTGTTCTGCC GCCTGTACTG CGCGTATTCC TCTTGCAAAA 11820
 CTTTGACAG ATGCACAAGA ACGCGCACTG CATCTTCCC GTGCTCCAGT CACCAAAATA 11880
 GAACCCCCACT CCACACAAAG CGTCGGGAAA ACTATCCGCG AGGCACCTGC CAATGCGCAC 11940
 CGCTGAGTAC AAACACGCAC CCTTCCTTTA CACCGGCTTA AGTGTGGAC AGAACAAACAG 12000
 TGTACTGTTG GCGCTGCTTG TTGCGCACGT GTTCGTCGTT GCAGCcATkc gCGACACGGT 12060
 CGcGCTTTTT TCCATCGTCA GTACCGAACT CGGCGCACTG AGCGCCGCGC TCGTTCAAAC 12120
 AcTACGCACA CCACATGTGC CCCTGAGCGA CTCTCTCGTA CTGGGCCTGC TCATCGGTGC 12180
 AGTACTCCCC GCACACAAAct CTTTTTGAA cACATTTGT GTCGCGTTCT GTGCCgTATT 12240
 TTTTACGCGC GTTTTGTTCG GTGGCAAAAT CGGGAAATTGG CTCAACCCCCA TAGCGCTTGC 12300
 CCCTGTCCCTC CTCCGCTGT GCACGGAGGG AACTTCCCTC CCAACGTCTG GGCGTGTCTC 12360
 TGTGTACAG GGAGCGATGT CTTATCCTCT TTTCTATTCT GCGCTTGTGAGTGGGACGC 12420
 CGCCGTGCGT ACGTGGTGCA ATACGCAGGT GTTCCAACCA CTTGGCTTTA CCCTCCCTGA 12480
 GGGAGCGTTG AGCGCCTGTG TGTTCACTCA GGCTGCAGCG CCTGGTTTC GCTATCCAGT 12540
 ACTTACCCCTT CTTGCTGCAC TGTGTGTATA CGCAgTGCGG GCGCGACGCT ACATCTGTT 12600
 GTGCGCGTTT CTTGTGGTGT ACAGCACACT GTTTTTTTa CCCGCACACG CACACCCCTGC 12660
 AmCCCTTGTT TCCCTCATAA AAAGCGCGC GCTGTTTACT GCATTCTTTG TACTCCCTGA 12720
 GCCAGATACG TCAATGCGCA CAAATGGCGG GGCTGGATC TCAGGGGAC TCTGTGcTAT 12780
 GTGCGCGTTT TTTCTTGCAA AAAAGAATAG TTCTGCCCA GATATGTGGG GTGCACACGA 12840
 CATGCACTTG TGGAGTGCGA TACTACTCAC CAACATCGTA CAGCCACTCA TTCTACGCGC 12900
 AGAACCTGG TACTACTATG TGCGGAGGCG TCGCTATGAC GTACAAACACT AACACGAGTC 12960

186

TTTCATCCTA CGCAGgATTG AGCGCATTG CGTTGTCAGT CTTTGCAATT CTATGGGCA 13020
 CCGCGCGCAC TGGTTCTTT TTAAAAGAAA AGGCGCTCAT CACyTGCGCC GCAGATATCC 13080
 TTGCAAGGCA AGCCCCAGAA CTTGGGGTCA CGTCACGCAC CCTGCGCATG GTACCGAGCT 13140
 CCCCCATACC GCAGgCTGAG GTGCTTCGGG GAAAAAAAGAA TACGGGAGAG GAAATATTCC 13200
 TATACTTTTT CCCACTCAGG GGAATGTACG GTTCGTTCC TACCCCTTTT TTGTACGATA 13260
 AAAAAGATGG TGCaCGCTTT TGCCaTCTCA TAGGTAATCA CCCTACACCG CGTGATGCAC 13320
 GCTTTTATGG CATATCGAgT tGCGCGCATC GCKyTTCAGT GTAGAAAAAT AGAACACCTC 13380
 CATCAAACAG TCGCATATGA GTAAGTACAC GGTTAACGCC GCGAGTGTAT TGTGCATTTC 13440
 TGGCATAGGA CTATTTGTTCTGCAACCCGG AACCTTGcc TGCGGTCTAC TACTCGTACT 13500
 TGGCTTTGG GTTCTATTTT TTCCTCGCT GCTGGCGAGA TTTCTCTCAC AGTTTTTAT 13560
 GCGCACGCCGC AGCgcTCCTT TGTCGAGGT CTGCTTACCC CTCAGCCA CCATTATGTA 13620
 TGACAACCTG ATCCAAGGCT TTTTCCCGCT TGTGCGTATG ATGCTGTGTC CTTACCTTT 13680
 CATTAmCsCG CTTTCGCGCA CACTCGATCT CTGCTTACCC GCATACGATG CAGATGCCGA 13740
 ATCGCTCGAA TGCGTAGGTG TCTTCGGCAT CATGATTGCG GGAATTCTC TTGTACGTGA 13800
 ATTAGTTGCC TTCGGGTGCG TTTCGCTACC GGCCCCGTCG GGGTCTTGC GCATCATCTC 13860
 TTTTCCACCC AGCAATGTAA TACGCTTGC AGCCACCGGC GCAGGGACCC TCATAAGCTG 13920
 TGGTATTGTT CTTGGATAT TCCGCAGTGC AGGTAACGAC CACACGCCCT CTTTAAGGAG 13980
 TGAATGGTGA CAATGGTGCC ACCCCCTGTTT TTCGTATGCG CCCCTTTCTT TGCGGAGGGC 14040
 ATCGGATTAG ATCGCCTGGT AnC 14063

(2) INFORMATION FOR SEQ ID NO: 2:

(i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 14244 base pairs
- (B) TYPE: nucleic acid
- (C) STRANDEDNESS: double
- (D) TOPOLOGY: linear

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 2:

CTCTGCTTGC CCCCTATGAG AAGACGGAGG CGCTTTCTCA CTCTTGCGC GTGTGCTGTG 60
 CACCTTCTTC CTCTTTCCC TCAGACGATT ACAACCGCTT CTGCTTTTT CGCCTCAGTT 120
 TCTGAGGATT ATGCCCGGTT AAAGGATTAC GCTGCTGATT TGGCCATGAG CACCGGGTCA 180
 GGAACCCGCG CGCACCTTAT GCGCGAAAG GTTATTTTA AATATCCAGA CCGTCTGCGT 240

187

TTGGATTCT CAAGCCCTGC TGAACAAACT ATTGTCTTCA CGGGAGATAG CCTGACCATC	300
TACTTGCCCA CCTCCCGCGT CGCGCTTGTA CAATCGGTAG CAAAAGATGA CACAGTAAGT	360
GCTGCTTCTC TAGCTCGCC TCATGGTCTT GCGCTTATGA AGCGGTTCTA CACGATAGCC	420
TACGAGACGA GTTCTCTCC TGTTCCCTG GGTCCGGACA GTGGGGAGAT GGCGTTGCA	480
CTGGTGCTCA ATCGTAAGTC TGCAAGCAGAA ACATTTAAAT CTGGCGCGT GCTTGTCTCG	540
GCACATACCA AGCTTATCCG TCGCATTGAA GCGTGGCCTC TTTGGGGGA AAAAATAACA	600
TTTGATTCTCA GCCACTATCG TTTGAACGTC GGTATTCCAG ACACACGGTT CCTCTACGAT	660
GTGCCCCCAA CCGCAAATGT GGTGCACAAT TTTCTCTTTG CTGATTGACC GCTGCCCCCA	720
AAGGACGTGA CGATGCCAGA TATTGGAGAG CTGCTAAAGA CGACCGCGA ACGCAAACAC	780
CTCAGTCTCG AACAGGCG CACGAGACGA GTATCGCACG CCGTTACCTG GAGGCGCTCG	840
AGAACGATGA GTATGATGTT TTTCGGCGC AACCTACAT CCTTGGCTTT TTGCGCAATT	900
ACTGCGAGTA CCTCCAGCTG GATACGGAGC AGTGCATCGC TCGCTATAAA CATTAAAAAA	960
TTCAAGAAAT GTCGCTGCCA ACGGAGACCC TCCTACCGAG TAAACGGTGG GGTCATTTC	1020
CCCTGTTAAA rGGAGTTGCC TGTGTGCTCT TCCTGGTGG GGTGCTGGGT GTGTATTACG	1080
CGCGGCACCG CnCnTnGGGT TTTCTATCCC GtATTGTGTT CTTTGGCAGA GCACAGCGTA	1140
CCCCAAGGGAA GCTGCTCTCCC CCCGATGCAA CGGGGGCGGT GCGCGAAACA GTGCGCTGT	1200
CTTCTGCACA ACATGAGGAG CGTGCACGAC GCACCGTATA CGAGCGCATC TCGCTATAAG	1260
CTTGCTGAGG AAAAGTTGAG ACACACGGTC TTTCCAGGAG ATGTGTTGGT TATCAGTTCC	1320
GGGGGGAAATG CGTACGAGCT AACCGTCAGC CGCACTACGC CGCACCTGTA TCTGGACACG	1380
CCCATTGGTA CACAGGTGAT CTCTCTTGGT CAGCGCTAG TGATGGATTT GAATACAGAT	1440
GTGCAGCCGG ACGTAGAAAT AAGTGTGGAA GACATTGAAG CACATCAGGC GGACGGGGC	1500
GCGCKTGTTC GCGTGTTCAG AGGTaGTCTG GTGCAGACGC TCCGTGATCG CAGTGCTCAG	1560
AGCTTTGTGC CTACAAGTGG GGTAAATGTC TCTGGTCAGA CGGGAGTCGC TGCCGGCGCG	1620
CGATATCAAG TTTTGTGAA AGGCAGGTGTT GCGTACCCGG TGACAATGAA CGCAACGTTT	1680
CGCTCGTACT GTTTGTCCG GTACGAAGCA GATCGCACGC GGCGGGAGGA CGGGTATTAC	1740
CAAAAGGGCG ACCAGCTGAC GGTGCAAGCA AACAAACGGGA TTCGGGTGTG GGCATCTAAC	1800
GGGAATGTGG TGCAGCTGCA AATTGTGCA GGCAGTAAGA CGGTGGATGT AGGCCTCAGC	1860
CGTCCGGGG AAGTGTGGT CAAAGACATC AAATGGATCA AAGATGAGGA CGCCGGGGCG	1920
TTCAAGTTCG TGGTCATGGA AGTAGACTAG CGCGCGCGGG CAGCAATCGC GTACGCKTTC	1980

188

CAGAGCGCGT GGACTGCAGT GCACAGTGCG CTTGCgCGCG CGCGGGAGCC GCTTCTTTTT	2040
TTCTCTCTTA CAAAAAGTAC CCGTA _g CGCT GCGCCCGCAG CTcCTGCAAA CAGCGTGG _c G	2100
CTGC _c TGCGG GCCGGTGTGC AAGAGCAAAG AGAAGGACTG ACAGTACCTC gCCACAGG _c G	2160
CGTGCAGTGC AGGAAGTGGC ATGGTGGCAC AGACGCTCAG GTATATAAGC GCGAAAGAGC	2220
ACTTCTTCGC TCAGAGCATT TAAAAAAAAGC CGTACATAAA AAGCTCCCCC TTCCCCCTCT	2280
GGGAAGGGAA ACGGTGAGGG AAAGAGAAAAA CAGAAGGAAG GAAATACTAG CGTGTGAGT	2340
ACGAACGCGT ATTCCGTAAC AGCCGCCGCA CGCGT _g CACT ACGTGTGGCG TGTTGCGCAA	2400
AGGGGAGAGG TGCATCAGCG TGGAACAGTG TAACAGAATC GGGCAGAGGG GGTACAGAGC	2460
GCATATATCC CCTGCAGTGG TGATGGCCAT TGCTGCGGTG GGAGGTTTTT ATGGGACTCA	2520
CGTGTATGAG GTACCGTTC CGTATGCATT CTPTTGTGCA GTACAGGCGT GTGTGCTGTG	2580
TATTGGGTGT TTGTTGGTCC GCAGTGGTGT GCGGTTCTTT TCTCGTTGGG GTGCTGTCCG	2640
TATCTGGAGG AGGTGGGGAA TCGCATACAC CAGCGTATGT CGGTGTTGTA ATACGCTTTT	2700
TTTCGTGTTTC TGTGGCTGTG TGTTGCCTG CGTTGCGCGA ACCTCCCTCA TGGTACAACA	2760
AGCTCCGTTG CAAACACTTG CACAACCCCA AAAACTACGC GTTTTGACTA TACACCTTTT	2820
GCAGAGGCCA AAGCCTGCAG GCACGCCGCTT TCGTGTTCGG GCGCGCGTAT TGGGTGCAGG	2880
TTACATAGAC GGTGCTTCCT TTTCTGCAGG TGGGGTGTGC ACTGTATTAT TTCCCTGCAGA	2940
GGTAATTGTTG CAGCACTACG CTACCGATAT GACGGACGAC gCGGATGCCCG CGCTCTGTCA	3000
GTATTACCGCG CGTGGGTTGCG GCTGTCAGAT TCGTGGCGC TTTGCATCTT CTGCACCGAA	3060
GCTTTTATC AGTAGTTCTA CACCACCAGC CTTTGTGTTGGC TGGAGTTCTT ATTGTCACA	3120
GATGCGCGCA CAGATGCGGG TTGCACTCAT GAGGTTTTA TCTCCATGGG GGCAGTGCAGG	3180
GGGATTGTTA CTCGCGCTCC TTTCTGCAGA TAGTGTGTTT CTTTCGGATG AAATGCGTGT	3240
CGCGTTTCGC CATGCAGGAC TTGCTCACGT GTTGGCACTC TCTGGCATGC ACTTGTCTTT	3300
GGTAGGGCG AGTGCAACGT TTTGGGCCG TTTCATCGGC ACAAGGCACA GAGGTATGCA	3360
GGGGGCGTTT TTGCGATGC TTGTCTTTGT GTGGTTGCA GGTATATCGC CTTCCCTTGC	3420
GCGTGCACCTT GGTATGACTT TAGTGTGTAT GGGAGGACAG ATGGCATACTG TGCGCGTAGG	3480
ACTTTTTCT GTACTGTGTG CTGTACTTAG CATACTATG CTCATTGCGC CGCATGATGT	3540
ACAGACGTTA AGTTTCATGT TGTCATACGG AGCGCTTGCA GGTATGTGT TGCTTGGCTC	3600
TGAGATTACT GAAATGATGT CGGGTTTGAT TCCTCGGCCA CTTGCATCGC TGCTTGGAAC	3660
GTCCTGTAGT GCGCAGTTTT TTACAGCACC GATAGTGCCTT TCGGTCATTG GATATTTG	3720

189

CCCCATTGGG GTACTTGCCT CGTGTGTGGT TAGTCGCCT ATGCCCTAT TTTTGATAGG	3780
GGGGAGCGTG GCGCTGTGCT GCTCTTGGC AGTGCCTGCT GTTGCCTT TTTTAAGTTG	3840
GGGTGTGTAC TTTTTGGTG AAGGACTCTG TGCGGTTGTG CGTTTTTTG CGTGTGCGCC	3900
GCTTGTGTAT GTACAGAGTG CCTGCAGACA TGTGTGTGCT GCATTATTTT CTTTTTTACT	3960
CGGTGGGGGA CTACTAGAGG CGGCAGCTCG CGTGCCTGTT CACAAGGATA CATATGTGTT	4020
GCCCGAATTA TAATTGGCG CGTGCACCTG CACAGTTTT GACGGAGCGC GGTTTGCGGA	4080
TGCATAAAA GTGGGGGCAG AATTTCTGC TCGATCCGGT GTTACGTACG CAGCTTGTAA	4140
AGATATTGGC GCCGGAGCGT GGGGAACGTG TATGGAAAT tGGTGCAGGC ATTGGTGCAGA	4200
TGACCGCACT TTTGGTGCAA AACAGTGATT TTTAACAGT GTTGAATT GATCGCGGCT	4260
TTGTGCAGAC ATTGCGAAA CTTTTGATG CACACGTCCG TGTGATAGAA GGGGATGTGT	4320
TGCAACAGTG GCATGCTGCA GCAGCACAGG AACAAACCTGC GTGTGTTCTA GGAAATTAC	4380
CCTACAATAT TGCTGCCGT TTTATTGGAA ACACGATCGA ATCAGGCTAT ATTTTTAACG	4440
GTATGGTGGT GACCCTCAA AAAGAAATCG GGTTGAGAAT GACTGCGCTC CCTGCACAAA	4500
AATGGTATTC ATACTTTCA GTACTCTGTC AGTGGCAGTA TGAAGTGCAGT GTGATTGCA	4560
ACGTTGCGCC TGTCTGTCTT TGGCCGCGTC CTCATGTAGT TTCTCAAGCA TTGGTACTCA	4620
CCAAGCGTAA TGCGGTGCCT TCTTGTGTGG ATCCTGCGCT TTTCTGCAC GTGACGAAAA	4680
CTTTGTTTC TGCGCGCGT AAAACGGTAA GAAATAATTy ACTCACGTGG CAAAAAAGGA	4740
TGCCAGGCGG TGCAAGCTGTG TGTGTTAGAAG AACTCTGCGC ACGTGCAGGT ATTGACGCGC	4800
GTGCGCKTGC AGAGCAACTG AGCATCTATG ATTTTATTAC GCTTTCTGAT aCgtGCGCG	4860
CGCTACTGTA GTCCGGTGTG GGTGTTGAAT GGCGCGTGT TATATTCTTT TTTTCAGTGT	4920
GTTTTTGTT TTCTCGCTCT TTCTGAAGA CGCCGCGCGC GATGTGAAAC CTAGCGATGC	4980
GCCTGTGCC C TATGAGGACA CAGAATTTTC CTTATGGCAG AAAGAATTGT ATCGTTTGAT	5040
AGCGCTGTCC ATCGGTGCAT TCCCGATAGT AACGCTGCTC TCTTTATCA CGTATGACAT	5100
CATACGTCTT ATTCAAGCAAT GGTCGACAAA GCCTCCGACA TGGTGGCGC TGATTATTCC	5160
TGGCGCGGAc TAcCGCCAcT GAGTACGAAG GAGCGCGCGA TAGTTTTGG TGTGGCAGTG	5220
GGGATTCTG TGACGATGG ATTAATTGAC GTGACGTATC GTGCAGTGAAC GCGTGAATA	5280
CACCGGGCGTA GTCTTGCAAC GTTGCAGTT AGTACCAAGAC CCGATAGAAC TGGTGCCACT	5340
TGATTCTTT TTGAGGGGA CTGACGATAG CACGTGAAGG TGACACAGGT GTCGTATTCT	5400
TGCAGTGCAG AAAGCACGTC GGTGTGCAGT GTTCATTTT CCGTTTTAG AATACCGGGC	5460

190

GCGACGTGTC GGTGTGATGT GTTCTGCGCA GAACATGTTT TTTTTTGTGC ACGATCTCAG	5520
CTCAAGAGGA TCGTGCCTGC ATTGTGTGTC AATGGGCATA CGGCAAAGTT TTCAAGACCT	5580
CTTCACGTGC GCGATCGGGT GTCTTTGAG TGGGTACGCT CAGTGCCCCC GGCGCTCATT	5640
CCTGAGAATA TATCGCTTTC TATTCTGTTT GAAAACGAAG ACATTATTGC GGTGAACAAA	5700
GCGCAGGGCA TGATAGTACA TCCTGGGCAGGCCACTGGAA CGGGAACACT TGTTCAGGGC	5760
CTCAGTTCT ACCGGGTGTA TCGTGCACGT TTTGAGGATG AGTTTCTCG TCAATTTCA	5820
AAAGGATTTT CCGATTTTTT CAGTACCCCTG CGTCAGGGTA TTGTGCACCG TTTGGATAAA	5880
GATACATCGG GCGTACTCCT CACTTCGCGC AACATGCATG CTCATGAGGC ACTTGTACGT	5940
TCGTTAAAAA AAAGACAAGT AAGAAAAGTA TATCTTGCCT TATTGCAGGG TGTTCCCTGCA	6000
CGCGGGTTG GGGTGATTGA ACAACAATC GTGCGAGATA GAAGACGACG CACCGGTTT	6060
GTTGCGTCTG AAGATTTTTC AAAAGGAAAG TACGCACGTA CGCGATACAA GGTGATGAAA	6120
ATATGTGGGG CGTGCCTTT TGTCCAGTTT CTATTGGATA CTGTCGTAC CCATCAGATA	6180
CGTGTGCACG CGCGATACCT AGGATGTCCC GTTGTAGGAG ATCCGTTGTA TGTTCCCGG	6240
AATATCTGTG GCATACCCAC AACACTCATG CTTCATGCGT ACCCACTACG GTTTGTTCTT	6300
CCGAGAACGA AAAAACGCAT AACGCTGGTA GCGCCCATAC CGCTTCGTTT TGTTGACTG	6360
ATACACCGAT TATCGGTTAG GTAGGGTGTG GCAGGTGCGT CGGTATATGC GTTTTACTTC	6420
AGCAGCTAAA TAGAAAGAAC CGATGGCAAG GAGCGTTTA CGTTGCGAA AACTGGCATA	6480
CAGTGCACGG GAAATAATTG ACGCAAAGTC TTCGCTCCAA AAAATTGGGA CTGTTGGTG	6540
AAAATGCGTA CGAACACGCGT GGTATGTTTT TTGTATATCT GCATGTTTAG ATGTGCCCGG	6600
TATGGTTAAA AAAATTTCGC TGGCGGCATG TGAAAAAAGA GGGGGAACT GGGACACTGC	6660
TTTGTCCGCC GCGCACGCAA AGAGTAAAAT GTATTGTGCA GAACGTAGTA AAGAAGAGAA	6720
CGTACGACAT GCGCACCGTA TACTCTGcGT AgTGTGCGCA CCGTCAATCA CTATGaGTGG	6780
ATCTTCCTGc ATAATTTCAA ACCGTGcTGG cACGTATGCA CGGGaCAGTC CCCGCTCGAT	6840
TAACGTTTCG CTCACGGTAG GAAATAAAATA TTTTGCCCGG CACGCAGCCA GTGCTGCATT	6900
TTTTGCCTGA ACAATATCGC ATAACGCGAG TGTGCAGTGT ATATTtCGAG CGAATAATCT	6960
GCCAACAGGA TGCAGGGCGT TAAACTGAG AGTTGCaGTG TGTGTGAAGT GTTTTATTGA	7020
ACTTTCAATA TGTGTGACCA TATCTGGTAA GTAAAAGAAG GGAGCATGTT TTTCTCGCGC	7080
GATATGTTA AAAACGTGCA ATGCATCTTC TGGCTGATCA AAACAAAAAA TAGGCGTATA	7140
GGGTTTGATA ATGCCGCCCT TTTCTTTGTC AATACTTTT ATACGTGTT CTAATATGCG	7200

191

CGTGTGTTCT TGTTCTATGG GGAGAAGGAG ACAGATACTA GGACAAATGA TGTTTGTTGC	7260
ATCTAGTCTT CCTCCAAGTC CTACTTCAAA AACGGACCAT TCCATGCGTT GTTGTGCAAA	7320
TAGCATGAAC GCCAGTAGCG TTATAAGCTC AAACCACGTC GCCTGGCCGT AGTCGCGCAG	7380
ATTCTCTGTT TTTTCACCG TGTGGTATAC GTGTGTGCAC GCGCTTGCAT ACTCGGCAGG	7440
TGAAAAAAAC ACACCCGCGC GTGTTATTCT CTCTCTCGGA TCCATAACGT GAGGAGAAC	7500
GTATAGCCCCG GTGTTGAATC CAATTTCAATT GAGTATCGCT GCAAgCATAAC GTCCGCTGGA	7560
ACCTTTTCCC TTCGTGCCCG CAACATGGAT GCTCTGATAT GCGTTGTGTG GATTACAAAG	7620
CGCGCGTGCA AGTGCAGTCA TCCTGTGCAG AGGTGGTGTG CCGCTTGGGG GCATTTTCTC	7680
AAGCGTGCAGA ATGCGCTCAA CCCAGGCGTA AAAATCTTGA AAAGAATGCA CCGGTATATG	7740
TGAAAATCCG TGTGGCCTCG GTCGCACTAT AATATGCAGT ACGGAAGAGG CAGCAATCCT	7800
TGCCGGGAAA GAGAACTGAT GTACATTGCT AAGGTCTGAC ATTTGAGCTA AAATCCGCC	7860
ATGAAGCCGG GGACTCTACC AAAAGATGTG TCAGGTATCA AGATTACAT GATTGGTATC	7920
AAGGGCACTG GCATGTCTGC GCTTGCAGAG CTACTGTGTG CACGGGGTGC CCGTGTGTCA	7980
GGTAGTGATG TTGCAGATGT GTTTTACACG GATAGGATTG TCGCCCGTTT GGGTGTCCC	8040
GTGCGTACTC CCTTTCTTG CCAGAACCTT GCTGACGCTC CCGATGTGGT TATCCACTCT	8100
GCAGCCTATG TGCCTGAAGA AAACGACGAG TTGGCAGAGG CGTACCGCG GGGTATTCCT	8160
ACCCCTACCT ACCCAGAAGC GCTGGGGGAC ATTTCTGTG CGCGGTTTC GTGTGGTATT	8220
GCAGGTGTTTC ATGGAAAGAC GACCACGACC GCGATGATTG CTCAAATGGT AAAGGAGCTG	8280
CGCCCTGATG CGTCCGTCCT TGTGGGGGAC CCGTGTGCGG GAAACAATGA TTCTTGTGTG	8340
GTTCTTAACG GAGATACCTT TTTTATCGCA GAAACGTGCG AGTACCGTCG GCATTTCTG	8400
CATTTTCATC CTCAAAAGAT TGTCTCACC AGTGTGAGC ACGATCACCA GGATTATTAC	8460
TCCTCGTACG AGGATATACT CGCGGCATAC TTTCATACA TAGATAGGCT TCCTCAATT	8520
GGTGAGTTAT TTTATTGCGT GGATGACCAAG GGCCTGCGGG AGGTAGTGCa GCTTGCCTT	8580
TTCAGTAGAC CGGACCTGGT GTATGTTCT TATGGGAAC GTGCCTGGGG CGATTATGGG	8640
GTCAGTATTC ACGGTGTTCA AGACCGGAAG ATAAGCTTCT CATTGCGGGG TTTTGCAGGT	8700
GAGTTTATG TTGCGCTCCC CGGTGAGCaT AGTGTGTTGA ATGCAACCGG TGCGCTCGCA	8760
TTAGCACTGA GTTTAGTGA GAAGCAGTAT GGAGAGGTTA CCGTTGAGCA CCTCACGCTC	8820
TGCggAAAGGT ACTCGCTCTT TTTCAAGGAT GCGGGGAAG GAGTGAAGTT CTTGGGGAAAG	8880
TGCCGGTAT TTTGTTCATG GACGATTATG GACATCATCC GACTGCAATT AAAAAGaCTC	8940

192

CGCGGGTTAA AAACGTTCTT TCCGGAAAGA AGAATTGTCG TCGATTTAT GTCCCATA	9000
TATTCGcgTA CCGCAGCCCT CCTCACCGAA TTTGCTGAGT CTTTCAGGA TGCGGATGTA	9060
GTTATTTGC ATGAGATTTA CGCCTCTGCT CGGGAAGTGT ATCAGGGCGA GGTGAACGGT	9120
GAACATCTT TTGAATTAAC TAAACGGAAG CACCGGCCGG TGTATTATTA CGAGGCTGTC	9180
ATGCAGGCAG TGCCTTTTT GCAGGCTGAA TTGAAAGAGG GCGACCTGTT CGTTACGCTC	9240
GGCGCTGGAG ACAATTGCAA ATTGGGTGAG GTGTTGTTCA ATTATTTAA AGAGGAGGTG	9300
TAAAGTCGG TTGCGGTTTC GCCAACATGG TGGTGGTGCC gGCTGTGGAT CTGGTGGATA	9360
TAGGGTGAAG TGAGACAGGC TCGGAATGGA TGGTGATGCA AAGAGCGGAG TGCGGAGGGG	9420
TGCGTGAGTA ATCGGTGCGA TGTGTCTGGA AATAAGGCGG TACgCATAGC AGTTTCAGGC	9480
GGTCAGGGT GTGGTAATAC CACCGTGTCT GCATTGCTTG CGGAAAGACT GGGACTTCCC	9540
CTAGTGAATT ATACGTTTAG GAATATTGCC CGGGAGTTGG GTATCTCTCT TAGTGAGGTG	9600
CTCGAGCGTG CGCGGACGGA TAATCATTG GATAAACAGC TTGATGCGCG GCAGCTCTGT	9660
CTTGCATGC GTTCTTCCTG CGTGGTAGGG TCGCGCCTGG CCATTTGGTT GGTGAAAGAT	9720
GCCGCGCTGA AGGTATATCT TTTGGCTTCA TTAAAAGAGC GGGTGAAACG TGTTCTCCAA	9780
AGGGAGGGAR GGGACGTACA GGATGTTGAG CGATTACGT CTATGCGTGA CGCTGAAGAT	9840
ATGACTCGCT ACAAAAAGTT GTATCGTATT GATAACACGA ATTACAGTT TGCGATCTT	9900
GTTCTAAACA CAGAAGGGTG CGATCAAGAA ACAGTGGTGA GTATTATTAT TGAAATGTTA	9960
CGCGCTAGAG GGATAGCTTG GTAGGGCTGA GCCAATCTGC GGGTGATATA GAAAAGTTTC	10020
AAAACGCCAT ATTGGATTTT TATGCACAGC AGGGCAGGGG TTTTCCGTGG AGAACTACTT	10080
GCGACCGCGTA TGnATACTGG TGTCTGAGTT TATGTTACAA CAGACACAGA CGGAGCGGGT	10140
GTGTCCGAAG TATGCGAAAT GGCTTCATCG TTTTCCCTCT TTGGAGTCTC TTGCGTGC	10200
TCCATTTGCG CACGTGCTCC AAGCGTGGAT TGGATTAGGA TACAACAGGC GCGCTCGTT	10260
TTTGATCAG TCGGCAAAAC TCATTGTTGA AAGGTATTGT GCAGTAGTTC CTGATGACCC	10320
GAGTGAECTA AAGAAGCTCC CCGGTGTCGG TGACTATACT GCCGCTGCAG TTGCTTGCTT	10380
TGCGTACAAT AAGGCCACCG TGTTTTAGA AACAAACATC CGTGCAGTGT TTATACGCTT	10440
TTTCTTTCCC GATACGCACC AGGTCAGTGA TCGGGAGTGTG CTCTCGCTGG TCCGGTGCAC	10500
CCTGTATGAG GAAAATCCTC GGCGTTGGTA CTACGCACTG ATGGATTATG GGGCAGTTCT	10560
AAAAAGGAAG ATTACAAATC CTAATCGTCG CAGCAAGCAT TACGTGAAGC AGTCACCGTT	10620
TGAAGGTTCT CTGAGGCAGG TGCCTGGAGC GGTTTTAAGA GAGATAAGCG GCATGCAACA	10680

193

CGCGGTGCGC GAGAAAACGC TTTtCGCAAA GCTGTCCtt GAGCACGAAA GATTGAGCCG	10740
CGCTCTAGAC TCGCTGGTAA GCGAGGGACT GGTAGTAAAA ACAGAGGCTG GGTATTCCAT	10800
CGCTGATTGA TTCTTTATGA CTCAAGACGC TTGAGTATT CACAAATAAA GATGCCCTTC	10860
TCTTTTATT CAATGGCATC TGATGTGAGG ATCAGCATGT TGGGCTGCTT TGCGTCAAGG	10920
CGAACGCGAT CGTTGGAGGT CTGTACCAGG CGCaGTATCT TTTTGAAGGC AGTGTGCAG	10980
ACCTTGTCAA ACTCAATCAA TAAGCTTCG TGTGTTCTT TGAGTGAGAG AATGGCGAGC	11040
TTTTCGCACC GCACCTTGAT TTCTGCCACG GTAAACAACC CAGCTGCTTC CTCaGGATA	11100
GGACCGAACC GGGTGATAGT TTCCGTGCGT ATGCGCTCAA GCTCCTCATG CGTATGAGCT	11160
GCAGCGATTT TTTTATACAG TTCCATTTTA ATTTCATCTG CGGCAATGTA CGTATGGGG	11220
ATGAACCCCTC GGTAATTAAAG ATCGATGACG GTTTCTATCC TTTGCTCGTT TGGAGCATGT	11280
TGGAGGCCTT CTATTGCCTC TTCTAACAGC TGTACATACA GGTCGAATCC GACTGAATAG	11340
ATATCTCCTG aTTGTTCTTT GCCTAATAGA TTTCCTACCC CGCGAATCTC CATATCTTT	11400
AAGGCGACTT TGAAACCCGC CCCAAGGTCA GTAAAGTCAG AGATCACCTG TAAACGTTT	11460
ATTGCAAGGT CTGAAAGTGC CACGTCGTGA TAGTACAGCA GATACGCATA TGCTTTTTG	11520
TCAGACCGAC CAACCGTCC CCTGAGTTGG TAGAGCTGGG AAACCCCGTA CATATCAGCT	11580
CTATCTATGA TGATAGTATT TGCATTGGGA ACGTCGATAC CATTTCAAT AATGGTGGTA	11640
GAAAGCAGGA GCTGGAACGT TTTTGATAA AACCTTCAA AAATGTCTTC CAGTTCTCT	11700
GACCCCATGA GACTGTGGC AACGCATATG GATAGCTCAG GCACGAGTTT TTGGAGCATA	11760
CACTTTACGG ATTCTAAGTT TTGATTCTG TTATGTAGGT AAAAATCTG CCCCTCACGA	11820
TCTAGCTCTT TTCTGATTGC AGTGGCAACA AGGTTGGAT CAAACTGCTG GATAACCGTT	11880
TCTATAGGTA GGCGGCCTTC AGGAGGGGTG GTGAGCAAGC TCATGTCTCT GATTTTGAGC	11940
ATACCCATGT GAAGCGTTCG GGGAAATGGGC GTTGCACTGA GGGAGAGACA ATCTACATTA	12000
GTTCATCT GCTTTAATTT TTCTTTATCC TGCACACCGA AACGTTGTT CTCATCGAGG	12060
ATCATCAACC CAAGATCCTT GAAGGACACG TCCTTTGGA TAAGCCGGTG GGTACCCACA	12120
ATAAGATCGA TATCTCCATG CGCGAGTTTG GCGAGTATGT CCTTTGTT AGATTTAGGA	12180
ACAAAGCGTG AGAGCTTCTC GATTCTGAGC GGAAAGTGT TAAACCGATT GCAGATTGTG	12240
CGAAAGTGTGTT GTTCCACTAG TAAGGTGGTA GGGGTGAGGA ACACCACTTG TTTTCCTCCC	12300
ATTACCGCCT TAAATGCCGC GCGCATTGCA ATCTCTGTT TTCCGTATCC GACATCTCCG	12360
CACACCAGCC GATCCATGGG GACGGCTTCT TGCATATCCT GTTGACTTC TTCAATGCAT	12420

194

ATGCGCTGAT CGTCGTTC TTCGTAGGG AATGCTGCTT CAAACGCATA CTGCCATTG	12480
TCATCTTTG GGAAGGCGTG GCCGCGCGTA GTTTTCGCA GAGAGTAGAG TTCCACTAGT	12540
TTTGCGCGA TGTTTCAAC AGATTTTTG ACACGTGCTT TTCTCGTTT CCATGACTTT	12600
GACCCAAGGC TATCTAAGTG AGGTTTGTTC CCTTCATTTC CAATGTAACG TTGCACCAGA	12660
TGTGCCTGCT CAATAGGGAT AAGGATCGTT TCTTCCTGTG CATAGAGGAG GTTTACGTAA	12720
TCACGTTCTG ACTGTGCTGT TTTTATGCGC TCTATTCCCT TAAATAAACC GATGCCGTAC	12780
TGCGCATGCA CCACGTAATC CCCGGGATTT AATTCCACAA ATGTGTCGAT AGGCGTGCTC	12840
CGTGCCTGTT GCACGTGATTG AGGAGTTTT CTGCGCGAC CGAAGATTTG GCCTTCTTGA	12900
ACGATCAGTA TTTTGAGAGC AGGAATGCTA AATCCTGCAG AAAGCGCGA AGGTAGCACA	12960
GTGACGTCGC AACCTTGAC TAGTGCTCTG ATGCGm-TGC CTGCTGCTCA CTTTCTGCAA	13020
AGACGAAAAC GTGCCATCCG TCTTTGAAA GACGGAGTAG CTCTTCTTGT AAGTAAGGAA	13080
TGTTACCGAA GAAGCTGCGT GCAGGATCGC TTGCCAAGCA TATACTTCG CACGCTGGCA	13140
GCTGTGGAAA AAAGTGAGTG AAATACACCG TGTGCAGGTG GAGCGCGCAG ACAGCGGAAA	13200
AATCGAGCAC TATGTGTTCT GGTTGAGGAT ACCAGCGCGC AGtACATGTT CGTGCCTGAG	13260
TTGCATTTTA TGGTAGAGGT TCCGACACTC GTCTTGGAGC GCGCGTGCAC CGTTGTGCTG	13320
GCGTTCTGAG TCAAGATAAA AGACGTTGG GGGTGAAGGG CTGTGGCGAA AATATTGAG	13380
AACGCAGgTG GGACGTTCAA AGCACAGTGG ATAGAACATT TCCTCCCCTT CATACTTTT	13440
TCTGTGGGTG AGTTCTCGA TACACGGGAC GCAGTGGCA GGACATTCGG ACAGTTTTG	13500
GAGATTTGG TGGAGGAACG CTATACGCTC CTCACTCCAA AGAATTCTT TTGCAGCGTA	13560
CAGTGTGCAC GCAGATAACCT CTTGCAGGAC GGCACACGTG GACACCGCCA GTATATGGAT	13620
ACGTTCTATG GTGTAAAAT CACACACGAT TCGGTACGCT TGTGTGTTGT CAGCAGCCTG	13680
CGCTGCGGCA GCGATATCGA GAATTCTCC CCGGAGAGAA AACTCTGCGC AAgcGcTGaC	13740
GTGGTCGACA CGTGCATATC CCCATTGCAT AAGCTGGCA GCGAGCGTGT GGATCTCGAT	13800
GTGCTCTCCC ACACGGAGG AGCGTTTGAG GGTACGCACA TAATCGAGGG GAGGAACGGG	13860
GGTGAGCAGT GCACGCTGGG TGAAAACGaA CrcGCATGgC gGtaTGCATC GcGCTGTGCG	13920
AGTGCACACA nnCTnCTnAC CCGGTGAGAG AACACGTGT CGTTAGGTGA GACAGGGCGG	13980
TAGGGCAGCG ACCCCCCACCA GGGCAGCAC GCGCGTAGGA ACTGCTGCAT GTGCAAGGTC	14040
GGTGCAGACG GCGGGCAGCT CCTGTGTTGG TAnGGACTAC GAGCACTATG TGTGCGCAAC	14100
ACGTGCGCAC GTTATTCGCC AAAAAAGTAG GACCGCAGTC CAACGGTGCA CCCTTTCAAA	14160

195

CGCGGTAGGG AAAAGCGTGC GGCACCAGCG AnnGCAGCAA TTGCTGGAAG CTCATTTCCA	14220
AGAAAATGGAG TATGCCACGC AACAA	14244

(2) INFORMATION FOR SEQ ID NO: 3:

- (i) SEQUENCE CHARACTERISTICS:
- (A) LENGTH: 2109 base pairs
 - (B) TYPE: nucleic acid
 - (C) STRANDEDNESS: double
 - (D) TOPOLOGY: linear

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 3:

AATCGACGAA AGGTATTGTA GCAGAGAGAA TACTCAAGCC ATGCGTAAGG AGAAAAGTAA	60
ATGGCAAGTT TAGATCTACC TAAGAGTCCC AATGTGTTTC ATCCCGAAAA GCCGAGTGCG	120
GTTGGGTCAA GGAATTCACT GGCGCAGGAC TGTCGTGACC AGCAGCAGGA GGTGAACCAG	180
CTAATAGAGG AAGAGACAAA CAAGATTCTG CACCACCTGA ACACTAAACT GCCGAAGaGG	240
TTCTCGAGCG TCTGGACGTA ATGGGTGGGT TGAAGGAAAA GTTGTATAAC TACTTCAACC	300
AGAATTACCA GAACATGTTCA ACCCGGTACA TGGTGACTGC GGAAGACGAA ATGCTGAAGA	360
AGGTCCGTGG TTTCATCGAC CGAGAGGAAA TGAAGGTGTT GAACCGTTAC ACGCCGAAGG	420
AGATTGCCAT CCTACTGGAT GAGGTTGCGG GAGCGGATAA GTTCAACACC GGAGAGATCG	480
AGAAAATCGAT GGTGAATATG TACGGGCACT TGCAGGGTCA TATAACAGCGG GGTGTGAATG	540
AGCTTGAGAC GCACACCAAT TCTTTGCTGC GTCAGAAAGGT TGATGTGGGT GCTTTTGTC	600
GC GGAGAGAA TGCGTATGCG GTAGTCAAGT GTGCGTCAA GGACAATCTT GCGCGTCTA	660
AGACCGTCAC TGACGTGAAG TTGTCTATCA ATATTCTGGA CTCAGAGTTA GTTAGCCCTA	720
TCTTCCATTA CCAGACGACG GTAGCGTACC TTATTAAGGA TCTCATCTCC AATCACTACA	780
TAGATGCCAT CGACAAAGAA ATTGATCGCG TGAAGGACGA GCTTATCGAC CAGGGTAAGG	840
AAGAGATGTC TGATAGCAGT ATCATCTTCG AAAAGATGAA GATGGTGAGC GATTTCACCG	900
ACGATGACTG CGAGAAmCCT GACAGCAAGC GCTACGAGCT TATTTCGCGG GAGTTGATGG	960
AAAGAATCAG CAATTTGCGC GCGGAAATTG ATCCGAAAC TTTCGACCAA TTGAATGTT	1020
GCGAGAATAT CAAAAAAATC GTTGACCTTG AGAACATAAG GAATCGTGGC TTTAACACGG	1080
CTATCAATTC GATTACATCT ATCCTTGATA CGTCGAGGAT GGGGTACCAAG TATATCGAGA	1140
ACTTCAAGAA TGCGCGCGAG CTTATCCTTC GTGAGTATGA TGACACAGAT ATTTCGAAC	1200
TTCCTGATGA GCGTTACCAAG TTGCGCTTAA AGTACCTCGA TAATGCTCAG TTGATTGAGG	1260

196

AGCGTAAGGG	GTATGAGGTG	ATGCTTCGTT	CTTTTGAGAC	GGAGGGTGGAT	CATCTATGGG	1320
ATGTGCTGCG	TACTAAGTAC	GATAAGTCTA	AGGCGTCTAG	GTTCATGGCG	AAGATTACCG	1380
ACTTTGATGA	CCTTGCTAAG	GTGTACAAGA	AGCATATAAA	GAAGCATTAC	AAGGATAAGA	1440
CTGGTGAGCC	CGTGTACGAG	GATATTGCGA	AGGTATGGGA	CGAGATTGCT	TTTGTGAAGC	1500
CTGCTGAGAC	CGAGGTGGAG	CGGATGAATC	GTACGTTGT	GTACGAGAAA	GACAAGATGC	1560
GAAGGAAGCT	TATTCTGATG	CGTGGGAAGT	TAAAGGGTAT	GTATGATTAC	CAGTATCCTA	1620
TTGAGCGTCG	GGTTATGGAG	GAGCGTCTCG	CGTTCTTGG	ATCCGAGTTT	AACCGTTTCG	1680
ATTACTTGGT	GAATCCTTTT	CACTTGCAGC	CGGGCTTA	GCTCGATATC	GACATCACGT	1740
CTATAAAGCG	CAAGAAGGCG	ACGCTCGACG	GTATGGCTAA	CGTGCTTAAT	GAGTTCTTGC	1800
ATGGTATCTC	TAAAGGATTT	GGGGACGCTG	CCTTTCGTTTC	GTTTAGTCGT	CGTCGTTCAA	1860
CGGTGCGTGC	TGATATCGGT	CAGAGTTTG	CTAGTGACGG	CAgTGCCGAC	CAGAAGGAGT	1920
CCAGCGGTAG	GGTGGCTTTT	ATGGATATGG	TAAATGAGAC	TCCTGCGCTT	GAGTCPTCCG	1980
TGGCCGCTGA	GCAGGGTGGAT	GTGCGCTCGG	ATGTTGGAAT	GAAGACGAGA	AAGGTGGCGC	2040
GGTGGATGCA	GGCAAGGGTC	GACGTGGTAG	ACGGTCTGCC	ATTCCGCAAt	CTAGCGAGAT	2100
TGTAGATAC						2109

(2) INFORMATION FOR SEQ ID NO: 4:

(i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 9848 base pairs
- (B) TYPE: nucleic acid
- (C) STRANDEDNESS: double
- (D) TOPOLOGY: linear

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 4:

CTGGACATGT	TTCCCTGCCTC	TGAACCTTGG	GTGAGGGAGT	TTGCACAGAG	GGTGGGGATT	60
CACGTGCAAG	AAGGTGCACG	GCTCGTGAAT	TTGCCTCGTC	ACCTAGCCA	AATCTATGAA	120
GCTTTTTTTG	AAGGAATTGT	GCTGTGGTGT	ATTTTGTGGT	GTGCGCGTCG	GGTAAAAACG	180
TATAACGGCT	TTTGGTGTG	TTTGTATGTG	GTGGGGTACG	GAGTGTTCG	TTTTTTTATT	240
GAGTATTTC	GTCAGCCTGA	TGCGCATTTG	GGGTACAGGT	TTTCCGCCAC	GCAATCGTCT	300
CCGATTTCACC	TTTCCAGTC	ATGGAGTGAT	GTTTCCACCG	GGCAGATTCT	GTGTGTTCTA	360
ATGATTCTCG	CAGGTTGGG	TGGGATGTT	GCACCTTCGG	CGTATCACAA	GGGGGATAGT	420
GTGCGGAAAG	CGCGTGTATG	AAAATGAAAA	GAATGCACCG	ACTGGTCCAT	CAGCCGAGAT	480

197

GGGGTGGCGC GATGTACCTT GCGTATAATG CGCAAAGGTG AGCCGATTcT TCTTCCGTG	540
TCCCTAATGAG CTGGTCTTCT TCTGGTGTAC CAAGCACCGC ATTGCCTGCC AGAGGGTTTC	600
GGTATGTGAT TTCTCGAAAG GGGTCAGAT AGTCTTGCA TCTCGTTTG TGAATTGTG	660
TTACTCTGTA CTTTCGATG GTAATTGCGT GTTGCACGTA GGTGCAAAGT AGCGGTGCGT	720
CTATCCTTGT TCGTGCAGA AACATAACA CTGCATAGGC GTGGAGTGCT TTATCCGTTT	780
CGTACGTGGC AGCAAGCGCG TCGTATTCTT TGTGGATGTG GTGCCAACTT TTGAACCTTC	840
CCGTTTCAAT GTTTGTAGA AGTGTCTCAA ACTTATCTTC AGGTACGAGT TGTCCCTCCA	900
TGTTTACCCA GTTGTGTGTG ATAGTTTAG GATCGTGCAGA AGTAAAGTGG GAGATGCAGC	960
GTGTTGTTG CTCAAAGAAA GACAGGAGCG TTTTGTGCA GTACCATATG AGTATGTCTC	1020
GGTATGCTTT GCGGGCTTCT AAGGGTTTTA GCACAAGTGT TGATCGGGTG GAGTGTCTA	1080
TGCCGGTGTGAG CAGTACGGGG ATTGTCGTGTG CTTGGTGGG ATGTTGAGG ACAATGTCTT	1140
CGGCAGTGAG CGCGCTGTTT CCCGCCTTAA CCCACGCAGC CTCTATCGCG CTGTCTAGTA	1200
AAGCGAGTGC GTTTCAATT TCTCCTATTG TGTCGGGAGC AAAGACAGAG ATTTCTACTG	1260
TTTGTGTTT CGTTTTCGT TTGTCGCG CAGAATT TTTCTGGCG TTCAATCGCA	1320
TACATGTTGT AGAGCCAGTA GTACGCGGGC ATGATTCTA ATCTGTTTC CCGTTCAATTG	1380
TTGGTGACAA GACTAAAAGG AAATGGAATA TCAAGTTCTG CGGGGTAATT TTTCCCCGTG	1440
ATTAATACGA AAGAACAAA ACGACAATTAG TGTTTGAGGG TACTTGCAAG TCCTGACCAG	1500
AACCCCTCGTC CCGCAATAAT TTCTCCGTCA TTTTGCGTG TATTGTGATT TGACCCAATG	1560
GTGGAGCCAG CAGCAATATT TGACTGACCA CGTATGAGTG CTGGATGAG GAAAGAGTTG	1620
TTGTGGTGTGTT GTTCATGGTA GGGAAAGATG AGTGCCTTAA AACCTCACA ACACGAAATC	1680
GTGGAGTTAT CACCTAAAAC TGAGTGAATG AGTCTGGTAC CATATTTAAG TGCGCAgTTA	1740
TTTCCCAGTA CAAAGCGCAC CGCCTTTACC CCATAGAACAC CACGGCAACC ATATCCGATC	1800
ACCCCGTTCA CTAACCTTAC TCCTTCTCCT ATTGCGTAG GTTCTGGAG AGATGACTGA	1860
ACAGTAAGGT TTTTAGTTT GTTGCTCCT TTTACGTATG ATCCAGGACC GAAGcACACA	1920
TCCTTGATGA TACGGCaGCT TTTGATAACG GATTGTGTTT CaATCGTTCC gTAGTATCCg	1980
CgGCGAGTGT CGTGGTGTG TTGAGTCATT GACTCGAAGC GTTGCATGAG CAATGTTCTG	2040
TCTCGGTGGC ATGCCACAA AAACGCGTCT GCTGCGATCA TCCCTACGAA GGGAAATATT	2100
TTCCGTCCGC CTGTTTCGTT GAGAGGATCA ATGGTGTGAG GGACTGCTTC TTGTTCTCCG	2160
TCTTTTATAA TTCCCTGCGCC GAATTTTGCGG TGGTTAGTGG TACACAGCTC ATCAATCCTG	2220

198

CTGAGGGATGA CATGATTTC C AATTATGTAG TGAGAAAATAT ATGCCAGTG ATGGATGGCG	2280
CAATTTTCTC CGACGTCGCA CGAAaTGAGT GTACTGTGGG TAATACCGGT TGGTACGGTA	2340
AAGTCGTGAT ATCGCAGAAA GgCTCGCTCG AGCGwasCGA TGCCTACGAG CCCTGCAAAT	2400
GATGAATTAC GTATGAGTGA CGCGTCGAAC GGATCTGCTA CTAAAACGTC GTGCCAGGTA	2460
TGCGCAGTGAT TGCCCTTTG TATAAGGGTG TGAATTTCTC CCTTAGACAA TGGTCTCCAT	2520
GCACGTGGGG GTTCCGGCCT CTGGGAAAAG CGGAGATAGT ACTCGTCTTT TCCCGGAGGG	2580
ATATGGGTTT GTGTGATGAA GTGGTATCCA AAAGAGGGTA GGTCTAAAAT TTGCACACGC	2640
ATTCTCCCTT TTGGATGCC ACTATAGTG CTGAATTTT ACATCTAAAT AAGGAATTGG	2700
GGTTGTGATG GGGATGGTGA TTTCCTGCAT GTTTACTTGA CATGACATAT TAGGAATGCC	2760
TAGATTGGGG CCCaGTCTTG TTTTTTAGCG TGCATTAGAA GTGGATGTaC TGGGGAGGAT	2820
CgTTGGCGGA TAACAAAAGC TTGCGGATT A TGGAAAGTAT TCGGGTACGA GAAGTGAGGT	2880
TGGTTGACGC TGTAGGGCAG CAGTGTkGGG TGGTGCCCAC CCCTGAGGCG CTGAGAATGG	2940
CACGGGATAT CAATCTTGAT TTAGTAGAGG TcgCTCCGCA GgCGAGTCCG CCGGTGTGCA	3000
AGATCCTGGA CTATGGGAAG TATCGCTTTG AGATGGCAA AAAGTTGCGT GACTCGAAAA	3060
AGCGACAGAG ATTGCAGACG CTCAGGAGG TGCCTATGCA ACCGAAGATC AACGACCATG	3120
ACATGGCGTT TAAGGCCAAG CATATACAGC GGTTTCTCGA TGAAGGGGAT AAGGTGAAAG	3180
TGACTATCCG CTTTCGTGGA AGGGAGCTTG CGCATACCGA TCTGGGTTTT AACGTGTTAC	3240
AGAATGTGCT TGGCCGTCTG GTGTGTGGGT ATAGTGTGA GAAGCAGGCA CCAATGGAAG	3300
GTCGGTCTAT GTCCATGACG CTCACTCCGA AGTCAAAGAA ATGATGGAGT GTCGGGTAAC	3360
TGCAGTTCGT GTTGTGGAT AAAGGGGAGA AAGTATATGG CTAAGATGAA AACGAAAAGC	3420
GCACAGCAA GCGTTTTAGT GTAACCAGGG CTGGTAAGGT AAAGTTCAAA AAGATGAACC	3480
TGCGTCACAT TTTGACGAAA AAGGCCCGA AACGCAAAG GAAATTACGT CATGCGGTT	3540
TTCTGTCAAA AGTTGAGCTT AAAGTGGTGA AGCGGAAGCT GTTGCCTTAC GCGTAGgTGG	3600
CAAGCGTGAG AGGACGGAGG AGCGTGGTAT GTCTCGATCG TTGAGTAGTA ACGGCAGAGT	3660
GCGCCGGAGA AAGAGGATTT TAAAGTTAGC CAAGGGCTTT CGGGGTAGGT GTGGCACGAA	3720
TTACAAGGCG GCGAAGGATG CGGTCTCGA GGCTCTTGC G CATAGCTATG TTGCGCGGAG	3780
GGATAGGAAG GGGAGTATGC GCAGtTGTGG ATCAGTCGCA TCAATGCATC GGTTCGTACG	3840
CAGGGtTGAG CTATTCGCG TTTATGAATG GTCTCTTGC G GCGTGGGATT GCGCTTAATC	3900
GCAAGGTTCT CTCCAATATG GCAATTGAGG ATCCAGGTGC GTTTCAGACG GTGATCGATG	3960

199

CTTCTAAGAA AGCTTTGGGG GGTGGAGCGT GCTAACCTC GGTCAGGTA AAGTGCTGGA 4020
 GGAGAAGGTT CGGAAGgCGG TGCACCTTGT CCAAATGTTG AAGGAAGAAA ATGCCGcgTT 4080
 GCGGgCTGAA ATTGATGGAC GTGGTAAGCG TATTACGGAG CTGGAGCAGC TGGTGCTTGS 4140
 CTTTCAGGAT GATCAGACGA AGATAGAGGA AGGAATTCTT AAGGCACTGA ACCACCTGAG 4200
 TACATTTGAG GATTCTGcGT ATGGAGAAGC GCTTACGCAA CACGCCGCGA AGgTTCTAGA 4260
 AAACCGGGAG CATGCCGGGC TGCTGAAGA ACTTACCAAGC CGTACCCAGA TGGAAATTIT 4320
 TTAGTGGTCA GTGTAAAGGG GCAGTTGCAC ATCGATCTGT TGGGAGCGTC TTTTCCATC 4380
 CAGGCTGACG AGGACTCCTC GTATCTCGT GTCTTGTATG AgCATTACAA GATGGTGGTG 4440
 TTGCaGGTGG AGAAGACGTC aGGGGTCCGC GATCCCTTAA AGGTcGCCGT GATTGCgGGT 4500
 GTGCTTCTCG CGGATGAACt GCATAAAAGAG AAGAGGAGAC GTCTTGTACA GTCCGAGGAA 4560
 GATCTGCTGG AAATAGGGG GTCTAnCCGA GCGTATGCTC GAATCCATCA GCAAAGTGGT 4620
 GGACGAGGGG TTTGTGTGCG GCGCGATTG AGGGTTGTGT CCTTCTTTGT GTACGGGACG 4680
 TCCTGCGGTG ACGCTGTGGG TGGACGGGA CTCATCCCCC GCGCgcGTCC GCGTACTTGT 4740
 CGCGAGAGCG GCAGCGCGCC TGGGGTGTGT GGCTCGATTT GTGGCCAACC GTCCTATCCC 4800
 TCTCGTGCAt AAAGCCGCATT GTATCATGGT CGAGACTCAA CCTGTTGACC AGGCTGCGGA 4860
 CCGTCACATg CATCGCGTAT GCGCGAGCGG GTGATTTGGT CGTCACGCGT GATATCGTGC 4920
 TTGCAAAGGC AATTGTAGAC GCGCGCATCT CTGTTATCAA CGACCGGGGT GATGTGTATA 4980
 CGGAGGAGAA CATA CGCGAG CGACTCTCGG TGCGTAACCT CATGTACGAC tGCGAGGGCA 5040
 GGGACTCGCC CCTGAAACAA CGTCACCGTT CGGCAGGAGG GATGCCGCAC GCTTCGCAGA 5100
 CTCCCTAGAT AGGGAAACCG CGAAAgcTCCT CGGGCTTGCC AGGCCGGGG AGGGGAAGAC 5160
 AGGGGAGGAG CAGTGCAGT GGCCCTCCGC GCAAGGGAAA AGCAAACCG GCCGCCGGTG 5220
 ACCGCACGCA AGACACTAAG AGTCCAAGGC CGGGCGGGTG GACTCCTAGT GTCTTCTACC 5280
 GCTTCTGCGA GATGAACCTTA AGCAAATCCA CCACACGGTT TGAGTAACCC CACTCGTTGT 5340
 CATAACAGGA CACTACCTTG AAGAAGCGCT TCTCGTTCGG GAGGTTGTTG TGCAAGCGTC 5400
 CCCTGCTGTC GTAGATGGAG GAGTACTGGT TGTGGATGAC GTCCGCGGAT ACAATATCCT 5460
 CGTCGCAATA CTGCAGGACA CCCCGCAGAT AGGACTCCGA CGCCTTCTTG AGCATCGCGT 5520
 TGAGGTCCGC AACGCTCGTC TCTTTTCCG TGCGGAAGGT TAGATCCACC ACGBAACCGG 5580
 TTGGTGTGCGG GACACGGAAG GCCATCCCCG TCAACTTACG TCTCGTAGAC GCCAGCACTT 5640
 CGCCTACCGC TTTCGCAGCT CCAGTGGTGG AAGGGATAAT GTTAACCGCT GCAGCGCGGC 5700

200

CTCCGCgCCA GTCCTTCAAA GAAACCCCAT CTACAGTTTT TTGCGTTCGG GTATAGGAGT	5760
GGATAGTCGT CATCAGTCCC GTTTCAATAC CGACTCCCTC TTTGAGAAAG ACGTGCACTA	5820
CCGGCGCGAG ACAGTTGGTA GTGCAGCTCG CGTTGGAGAC GACCTTGTGC TCAGCAGGAT	5880
CGAACTCATG CTCGTTCACCC CCCATTACAA TAGTCTTCAC CGGCTTAGAC GCATCCGAGC	5940
TCTTAGCCGG AGCACTGATG ATGACTCGCT TTGCTCCTGC TTCAAGGTGA CCGTATGAAG	6000
ACTCATTGCG GTAAATGCCG GTGGscTCAA TAACCACCTC AATACCAAGA TCCCTCCAGG	6060
GaAGTTGGGA AGGCTTTAAG CCGCGACCGC AGACACACTT GATCGATGC CCGCCCACCT	6120
CGAGGATATC CTCGGCAGGA GCACTGAGAC TAGAACCCAT TTGCCCCGTACGGAGTCAT	6180
ACTTTAGCTG ATAGGCAAAG TAGCGCGCAT CGGTGGAAAG GTCTACAACGCCG	6240
CGAACTCTTT CCCAACAGC TTCTGtCCGC CATGGCCTGG AGTACGAGAC GCCCGATACG	6300
CCCAAAACCA TTGATTGCAA CTCTCATTTC CCCAACCTCC TCTAAAAAGA GCACACATCC	6360
CGCGCAACGC TATCTGAAAA AAGATCGGCA CGTCAATCCC TCTTGCTGT AGGGCTCCCT	6420
TGCATTTTC TATGTGCCA GATACCATGG CCTCGCCTTG GAAGGTCTGG CCTCTAGTGG	6480
AAGATTATTA CCGCGTGCTT GGTGTGTCGC ACCGTGCCCTC GACCCCTGAA ATTAAGTGTG	6540
CCTTCAGAAA GAAGGCAAAG GCGTTACATC CGGATCTCGT TTCCCATACT GCAGAACTTG	6600
AGTGCAGGGC GGTAgCgCGC GAGCGCgCTC TTGCGCGTAT ACTCACCGCA TACGAGGTGC	6660
TCTCTGATCC GGGCGTCGC GCGAAATTG ACCTCCTCTA CGCGCGTTTC TGCGCACGTC	6720
CTGCTCCAGC GGGCTTTGAC TACCGCGTGT AmCTGCGTGC GCAGGTACGC TCTGCGCGAT	6780
GGTGGAGCTT ATCTTGTGTTG ATCTCTTICA CGGTTTGAG TGTGACGCTG TCCGGCGTA	6840
CTTGTCCCTC AAGTGTGGC CAGAAGGGTT CAACCTCGCC ACTCACCTTA CACGAGAGGA	6900
TTTTATGGAC TGTGGCTTG TGCTCGAGA GGAATTGCAT GTACGGGGAG AGTGTATGA	6960
ATGCTTTACT TTGCTCCAGG ACATCGTTTT TGAAGAATTG CGGTGCGCGT ATTTTCGTCA	7020
TTTTTTCTT GAAGTACTGA AGCTCGCTGA GCATATCGCG CTCGGTAcTG CGTCTGTGCG	7080
TGGTCGCAAC GGTAAATCCT GCGTATACTG CGCGCGCGCC ATGCCCTGCTT GCCTGCGCAA	7140
GAAATTGTCA CCTTCTACGC GTGCTTAGTT GAGTATTACG AACGTACGGG AGACCGCAGC	7200
GTGCGCGTGG CTATGCCAG AAGATGGATT CTGTCAGGTG AATGTTTGAC TGACCCCTGG	7260
CGGAGGGAGTA CCGTGGCTCT GGGGGACCTC CGAAGGCTGG AGGTCCCCCT GCAGCTAGTG	7320
AACGGACAGA GGAGGGACGC TTGAGCAGGA AGGAAAGGAC CTCATGATCC GCATTAAC	7380
ACCAGAACAA ATCGACGGTA TCCGTGCCTC TTGCAAGGCA TTGGCGCGCC TTTTCGACGT	7440

201

TCTTATTCCG CTTGTCAAAAC CGGGCGTTCA AACCCAGGAG CTTGATGCGT TTTGCCAACG	7500
CTTCATCCGC TCAGTCGGTG GTGTTCTGC CTGGTTCTCG GAAGGTTTTC CTGCCGCTGC	7560
TTGCATTCA ATCACACGAAG AGGTCATCCA TGTTTACCT TCAGCGCGTG TGATTCAAGGA	7620
CGGGGATCTT GTTTCCCTTG ATGTTGGTAT CAACCTCAAT GGATACATTT CTGACCGCGTG	7680
TCGTACTGTT CCTGTCGGTG GAGTTGCACA CGAGCGACTA GAACTTTGC GTGTAACCAC	7740
TGAGTGCCTC CGTGCAGGCA TTAAAGCGTG CCGTGCCgGA gCnCGyGCGC GCtgTTTCTC	7800
GCGCTGTATA CGCTGTTGCA GCACGGCACC GCTTGGCGT GGTGTACGAA TATTGCGGAC	7860
ATGGCGTGGG GCTTGCCTG CATGAGGAGC CGAACATCCC CAATGTGCCT GGCTTGGAAAG	7920
GGCCTAATCC ACGTTTTTG CCCGGTATGG TAGTCGCGAT AGAACCCATG TTGACGCTTG	7980
GCACAGACGA GGTGCGCACC AGTGCAGATG GCTGGACGGT GGTAACGGCA GACGGATCGT	8040
GTGCCTGCCA TGTGGAGCAC ACTGTGGCAG TTTTGAGA CCACACGGAG GTTTAACAG	8100
AACTACGGAA GTAGAGCGTA CCGGCTAGTC AGTATCTTA AGTGTGCGCG GTGTGCTGAT	8160
AGTACATGCA GGGAGCAGTT TGTGCACGGT AGGCAGCGTG TAAGTGTACG TGGCGGGCAC	8220
AGGTGAAGAG GGGATAAACT CGTAACCATA TCGCTGTGTG CTGCTTTAA CCCGGCTGT	8280
GTCGGTAGGG GTTTGGGTAC GCGCAGGGAC GTGGAGGGAC TCATGAACAT ATTGTTACC	8340
TCGTTGTGT GTGGGGTACA TCGGGTATGC CGCAGTTTT TTACAGCAGC GGCGTTGCTC	8400
GTTTTATCT GCTGCTCTGG TCATCCAAGT TCTGCGCGTG TGCCCTCTGC AGACACGATA	8460
GCTCGGGCG TTGCCGGAGA CAGTGGGAAC gCTGGGGGC GGACATTACT TCCGTGGGG	8520
GTTTCgCGTG AATCGGTGCA GCTGTTAGAA CGGCTGCAA ACACGAAACCG TCAGGTAACT	8580
GCCGAAGTGC TGCCTTCAGT AGTGACGCTG GATGTGGTGG AGACCAAGAAA GGTTGGGTA	8640
CGTGATCCGT TTGGCGGTTT TCCGTGGTTT TTCTTCTGTG GTCCGTGAAGG TCCGGGTGCG	8700
GGGnCTGGCG GTGGTTCTGG AAACAAAGGG GAAGCTGAGG AACGGGAGTA CAAAACGGAG	8760
GGACTTGGTT CTGGAGTCAT TGTAAAGAAG ACAGGGAAAGA CGCATTACGT GCTTACCAAC	8820
TATCACGTGG CGGGTAAGGC TAATGAGATA GAGATTAAC TGCACGATGG CAGAACCGTA	8880
AAAGGTAAAC TTGTCGGTGG TGACCAGCGC AAGGACATCG CGCTGGTCTC CTTTGAGGAC	8940
GCAGACCCAA ATATCCGTGT TGCCGTCTT GGTGACTCGG ATGCAGTACG GGTAGGAGAC	9000
ATTGTGTTCG CAGTGGCTC TCCTCTTGGG TACACTTCCA CTGTAACGCA GGGGATTATC	9060
AGTGCCTGG GTGCTTGGG GGGACGGGC AACAAATATTA ATGATTTAT TCAAACAGAT	9120
GCGGCCATAA ACCAGGGCAA TTCCGGGGGA CCAATGGTCA ATATTTATGG CGAAGTGATT	9180

202

GGGATTAACG CGTGGATTGC CTCCCTCAAGT GGGGGATCGC AAGGGATTGG TTTTTCAATT	9240
CCTATCAATA ATGTGAAGTC GGATATCGAA TCATTTATCC AGTACGGGCA GGTGAAGTAC	9300
GGGTGGTTAG GCGTGCAGCT GGTGGCAACG GATGCGGACA CCGTAGCATC GCTTGGTATT	9360
GCAAAGGGTA CAAAAGGGT GCTTGGCGCG GAAATTCTCT TAGTTCTCC TGCGCACAAG	9420
GGGGGACTGA AACCGGGCGA TTACTGTGTA AAACGTAAACG GAAAAGAAGT AAAGGATGTA	9480
AATCAGTTG TGCGGGATGT CGCGCGCTG CGCATTGGC AAACAGCAGT ATTGATTTA	9540
ATTCGCGGTG GTGTGCCGAT GACGCTTCG GTGCGCATTA CGGAGCGTGA TGAAAAAATA	9600
GTAAATGACT ACTCAAAGCT TTGGCCTGGG TTCATCCCAC TGCCGCTTAC GGAGGCCGTG	9660
CGTAAACGTT TGGATTTGAA AGCGTCGGTG CGTGGTGTGC TAGTTAGCAA CGCGCAGAGC	9720
AAAAGCCCTG nCGGGCGCTGA TGGGATTGAA GTCGGGGGAC ATAGTAGTGG CGGTCAATGA	9780
TCAAAGAGTC TCGAGCGTGC GTGAGTTTA CGCGGnGCTT GCACGTCAGA CGAGGGAAAGG	9840
TGTGGnTT	9848

(2) INFORMATION FOR SEQ ID NO: 5:

(i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 7415 base pairs
- (B) TYPE: nucleic acid
- (C) STRANDEDNESS: double
- (D) TOPOLOGY: linear

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 5:

CAAAAGGAGA TGTCAATTCTA TGTGGAGCAT GCCATTGCGG TATGCACCTG GGTTGTGCGC	60
AGACACTTCC GCCTGATGTG CGCACCAAGAT AGACTATGTT GTGGAACGCA CTCGGTCTCA	120
TCGGGGAACT ACTGTTGCGC TTGCTATAAA TTATGGGGA AAAGATGAAA TTTTACGTGC	180
GGTAAAAAAAG GTTTTGTGCA GCACTTCGTG CCCGGATGGT GAGCTTCTCA CCGAAGAAC	240
TTTCGGCGCG TGCCTTGATG CGCCGCAGTT GCGAGTGTC GACTTTCTCA TCAGAACAGG	300
GGGTCAGCAA CGCATGAGTA ATTTTTTGCT TTGGCAAAGC GCGTACGcGG AGTTCTATTT	360
TACCGATATC CTGTGGCCTG ACTTTGGGT AGAAGACATG cTGGCGGCC TGGATGAGTA	420
TCGCCTGCGC ACGCGTACCT TTGGGGGTTT GGAATGAGCG CGGAAATAAA GAGGCTGTTA	480
ATCTTTTTTT TCGGCGtTCC AACTATTCTT ATGTTGGTAT ATGCGGCACC GCATGcACAC	540
TTCCTAGCGT TCCATTGCT TATCTTCGGA TCAGTTATGG GTGCCGTATG GGAAATGCAT	600
GCGATGGTGT CgcGcAGGAT GTGCACGTAC CCACGGTTT TGTGATCCC TTTCAGTCTT	660

203

GTGCTTCGGC TTTTAGGATA TGCAAGCGCTG TGGCAGCCTG CACGGGGCGC TGAATCTGTC 720
 CTTTTTATTG GAGCACTGGG CACGCTGCTC ATGAGTGT TTTCACCGA ATTGGTGTAT 780
 TCGTTTCTG CTTCTTTGA AAACGCCCTT GAGCGTATGG CCTCGGCACT GTTGCTTGTT 840
 TTGTATCCAG GTATCTTAG CCTTTTTTTT TCGCTCATTA CGCGTGGCG TCATGCAGAG 900
 ATCGCaTTGG TAATTTTTT yCTCATGGTT TTTACGTGCG ACTCTTGTGC ATGGTTCTGT 960
 GGGACGCTCT GGGGAGTCAA CAACAGAGGG ATAATTCCCTG CAAtCCTAAA AAGAGTATgC 1020
 AGtTTTATgG AGGTTTTgCC GGTTGGTAG GTGCAGGgTG TTTTGGCTCA CTtGTATTG 1080
 GTTCGCgTGT GACGCTCTCT TTGGGGATGC TCATGGGTGT TGGAGCCTTG GTAGGACTGA 1140
 CTGCCATTGT AGGCATCTA GTCGAGTCGG TGATGAAACG TTCGGCTCAG GTAAAGGATT 1200
 CAGGATTTTT TACCCCCGGG CGGGGCGGAA TTATGGATAA CCTGGATTG tTGCGCCGTC 1260
 ACTGGGGACT TTTTACATTG CATGTGAGTG TTTTGGGATC GCTGCAGTAT GAGTGTGCGA 1320
 CGTGTGGTAG TGCTGGCAT TACTGGTTCT ATTGGAGCTG CAGCACTCAA ACTTCTGCGT 1380
 CGGTTTCCCG ATCGGTTCTT GCTGGTGGGC GCTTCAGGTC ACCGGCAGAC CGAGTACGCG 1440
 CGGGCGTTGG CGCGCGAGTT CTCTTATCA GATATCACTA TGACTGGCTC ATGTTCTGAG 1500
 CAAGAAGGTC GCGCACGCAT AAAGCGTCTG CTTTCTTCCT GTGAAGCAGA GGTGGTGGTA 1560
 AACGGTATTG CCGGCCTG TGTCCTTTT GCCTCTCTTG AGGTGCTCAA GACGCCTTGT 1620
 ACGCTCGCGT TAGCAAATAA AGAAAGTGTG GTACTTGAG CTTCTCTTTT GCATGCTGCG 1680
 GCACGCGAAA GTGGGGCAAC AATCGTTCCCT GTAGATTGAG AGCATGCTGC TATTTTCAA 1740
 CTTATTGAG cGCACGGCGC GCATGCGGTG GCGCAGGTAG TGCTCActGc GTCAAGGTGGT 1800
 CCATTTAGAA CCTTTCAAA GGAGTGCTTA GCGCATGTCA CGGTGGAAGA TGCGCTTCAA 1860
 CATCCGACGT GGCgtATGGG GAAGAAGATT TCTGTTGATT CTGCAACACT TGCAAATAAG 1920
 GCACTGGAAG TTATAGAAC AGTGCAGTTT TTTCGTATAC CGGTGGATCG GGTCaCGGTG 1980
 GTGGTGCacc CTCAGAGCaT AGTGCATGCG CTGGTGcAAT GTCATTGGG AGAAACGTAT 2040
 GCGCAGCTTT CTGTCCCTGA TATGGCGTCG CCGTTACTGT ATGCGTTGCT GTACCCTGAT 2100
 GCGCCTCTGC GTATCAAACCT CCGCTTGATT TTACATCGGG ACTGTCTTTG CATTGAAAC 2160
 CTCCGAGGGT AGATGACTTT CCGCTGTTGC GTATGGGTTT TGATGTTGCA CGGGCGCAGC 2220
 GTGCGTATCC TATTGCTTT AATGCAGCAA ATGAGGAGGC GGTGCGTGCg TTCTTGCAA 2280
 GAAACATTGG GTTTTAGAT ATCGCACACG TGACTGCACA GGCGTTGCAA GAAGATTGGC 2340
 GCGCAATTCC CAAACGTTT GAAGAAGTTA TGGCgTGCGA TAmGCGTGCg CGGATGTGTG 2400

204

CgCGGACGTG CATTGCACAG AGGTGGAGAG AGAGGTGATT AAGATAATTA TTGGCGTTGT	2460
GGTGCTTGGT ATTGTGGTGT TGTTTCATGA ACTGGGGCAT TTTGTCGCCG CGCTTTGGTG	2520
TCGAGTGGAG GTGCTCAGTT TTTCTGTCGG TATGGGGCCG GTCCTGTTTC GAAAGAAATT	2580
TGGAAAAACG GAATATCGCC TTTCGATGCT TCCTCTGGG GGGTATTGCG GTATGAAGGG	2640
AGAGCAAGCG TTTCAAACGG CGCTTGATCA AAAACTTCC CGTATTCCCG TTGAGCCCCG	2700
TTCACTGTAT GCaGTAGGAC CGCTCAAACG CATGGGTATT GCCTTGCAG GACCCTGGC	2760
GAATGTGCTT ATGGCGGTAA TGGTATTGGC ATTGGTTAGT GCGCTTGGCT CGCGTGTACA	2820
CACATTTGGA AACCGTATTG CACCGGTGTA TGTATACGAT AGTCTGATA ACTCGCCTGC	2880
ACGCCGCGTG GGACTTCAGG ACGGGGATAC AATcCTGCGC ATTGGTGACC AGCCGATACG	2940
CTATTTCACT GATATTCAAA AAATTGTATC ACAGCATGCG CAGCGTGCAT TGCCATTG	3000
GATCGAACGG AGGGGGCAGC TTATGCACGT GACCATTACG CCTGATAGAG ATGCGCATA	3060
TGGCATGGGG AGGGTTGGTA TTTACCATTA CGTACCGCTA GTTGTGCGG CGGTTGATGC	3120
ACACGGTGCT GCATCGCGGG CAGGTCTTGA ACCTGAAGAT AAAATTCTTG CAGTAGCAGG	3180
ACGCCGTGTG CAACACcAGT ACAGCTCCTT GCGCTGCTCA AGGAATTTCG AAAAAAGTCA	3240
GTCGTATTGA CTGTGCTGCG TTCAGGAAAG AGGCGATATC ATACCATTGC GTTGTGCGC	3300
ACAGAAAACG GGGCAATAGA TGTGGTATC GAATGGAAAG CTCACACCGT GGTTATACCG	3360
GGAACTTCTT TTTTGCAAG TGTCCGTGCG GGCATTGCAG AACCGTTGCG TATGTGTGTA	3420
TTGACGGTGA AGGGTATTGG TATGCTCTT CGGGGCCTGC AATPTCAGCA GGCTATCTCA	3480
GGCCCATTAA GGATTACGCA TGTGATAGGA GATGTGCC AGCATGGTT TCAGGAGAGT	3540
TTTTAACGG GACTGTCACA ATTATGCGAG TTTGTGGCAC TCGTGTGCGT CTCTCTCTT	3600
ATTATGAATC TACTCCCCAT TCCGATCCTG GACGGGGTT TGATTTTATT CGCATGTGTT	3660
GAATTGTTA TGCAAAGAAG CATACACCCG CGTGTGTTGT ACTATCTGCA GTTTGTAGGT	3720
TTTGCCTTGTG TTGCATTGAT ATTGTTATGT GCGTTTGGAA ACGACGTGAA TTTTTGTTT	3780
CACTAGGAGT GAGTGATGCA GTTACGGTGT GCGTGTGAGC GGGTGTGCA TATTGAACAT	3840
GAGACGGTAA TTTCCCTTGA TGAGCACCCG GAATTGTTG CGCGTATACA GCAGGGGGAT	3900
TTTTAACGTT ACCAGTGTCC GGCATGTGGT GCGCGTATTG TGCCGAAAT AAAAACAGAA	3960
TTTGTGTGGC ATGCGAAGAA TGTGCATTG CTTTGTGGTC CTGaGCGAGA GCGTTTGGCG	4020
TGTTTGGCTT TTTGTGCCGG TATGCATATG AGCGACGGAG ATAGTGCTGA CTTTGTGAA	4080
CCCTTTGTCT TACGGGAGCA CCAGACACCC GTGATTGGCT ACGCAGAACT TGCTGATCGT	4140

205

GTTGCAATAC TAGCATGGGA TTTGAACCCCT GAAATTGTTG AAGCAGTAAA GTTTTTGTG	4200
TTGGAAGGGG CACCGCATCT AGGAGACAAG AGAGTTTCGT GTTTTTTGTA ACGTTGTGTC	4260
GGGGACACCG GATCGCGCGT GATGGAGTTG CACGTGTACG GTATCAGAGA ACAACAAACG	4320
GCAATTATGC CGGTTCCCCT GAATGTGTAT GAACGCGTTG AGCGAGAGC mGGTAAACAA	4380
GCGGAGTTGT TTGAGGCCT GTATGTTGGG GCGTATCTT CATACAAGAA TGTTTTACT	4440
GACCGCGTAGC SCCGCACAGC GAGCAGCAGC TGGTGTGCGT GGTGTATGGG GTGTTGTACG	4500
CTTGGGCCTT TTGCAGACAG TGTAGAGAAG CGCGCAGcgA AGGATGTGTT TACTGAACCG	4560
GCGCGCTTTT ATCCCTCACA AAAATCAACG CTTGAATCTG CCCGGTCTGA TACATCTGAA	4620
TCTGAGAATG CATCTCTTC CGTTCCCTCC CACAGTCAGC AGGAGTTGGC GCCAGACTCT	4680
GCCGCGCCTG CGCGTAACTC TGTGTTGTCC CCTGCTCCTC CTGAAAGGAG AGAGAACGAG	4740
GGGACTGCGG TGCATGGGGC GGAAGTGACG CGGGCGGGAG CTGTCAGCCC GCGTTTTGTA	4800
GGGGGGCTGA CAAAAATACT GGCCGCCTCT GACCATAACAT TCTTCGCTGC AGGAAATGAT	4860
GGGTTCTCA CCCAGTACAC GTATCCGGAT TATAAACCGG ATACGTGGCA GATCACCCCT	4920
GTTTCTATCA AACACTGTGC AGTGCATCCG GACCGCGCGC GTATTGCCGT ATATGAAACA	4980
GATGGACGCA ATTACCAACCG AGTCAGTGTG TGGAAATTGGC GCACGAAAGA AATACTTTTT	5040
GCAAAGCGTT TTACCGCATC GGTTGTGTCA CTCTCGTGGG TTGTGCAGGG AAGTTTTTTG	5100
AGTGTGGGAA CAGCATCGCG CGAAGgTGTG ACGGTGTTAG ATGGGAGTGG AAATACTAGTT	5160
TCTCTATTTT CGGAAGAGCC TGGGGTGGTG TTGTTGACTG CGAGTGGACC GCGCCTGTG	5220
CTCAGTTATG CAGAATCTGG ACGCCCTCAGC TACGTAGATT ACAGCAAAAA GACAACCGTC	5280
AAACGTCTTC TTACCGAAAA GAATCTCCTG TCTCCCATGT TAATACATAA CGGTGCACAT	5340
CTTGTGCGTT ATAGAGACCA ACGTGTGTAT GTCATCCAGT CTTCAAGTGG CGCGGTGCTC	5400
ACCGAGTACC CTGCACGGAG tGcATGttTTT GCGCATAACAT TCAGCGATAG TCTTCCTGTG	5460
TGGATAGAGC CTGCTGAGTT GAAGTATCAC TGGCGTATAC GGAAAGcTGC GCAGCGTTCT	5520
GCTGATTTTA TGCTTCCTGA CAATGCTCGC ATAACAAGTG CGTGCTCGGT TCGCACGCGG	5580
GTCATCGTAG GAACCGATCG CGGGATCCTC TATGAATTGC AGCAGGGAGA TGACAGGGCG	5640
GTAACATATCC GCGCACTCAA TGGCGAGCGT CAGATATACG CAAGCGATGT ACATGGTGCA	5700
GATGAGGGCG CGTATTTTTT AGCAGACGGA TCCCTATATC ACAGCATGGC GTCCGGGGGA	5760
CCGTATCGTG TTTTGGTGCG CGGAGTAAAA GGAACTCGGT TTCTGCCTTA TCGTGATGGT	5820
TTTATTGTGT GGTCTGCAGG GAAAGAAACA GAGTTCTTC ATTGTGCGCA AAAGACGAGT	5880

206

CAACACAGGA TGATATATCG CGCGCGTTCC ACGGTAAGCG GCGTGTGGT GTATGGGCGT	5940
ATGTTGGTGA TTACTGAACC TTTCTCTGGA GTATCGGTGG TGGATATTGA GCGGGGGATA	6000
CgAGTTTTTT TTCACAAAGC GATTGGTATG CAGGATTTCGC TATTGATTAC TGATGACGTA	6060
ATTGTAGCCA CTCAAAGCGG TTTGCAGCCA CTTGTCCTGC TGCAATATGCC TACGGGGGAG	6120
ACATATACGC AGCGGTGGGA CGCGATTTCGC CTTGGCGTCC GCGCCATGA TACACAGCAT	6180
GTATATTTTT TTTCGTTGGA TACGAATGCC GGCACGACTG ATTTGATCCA TTTCGTCTGC	6240
AACTGCAGCA ACCCACAGAA AGTGTGTGTC GACGCATCCT CTCTTATAAG GATGAGGATA	6300
TAGATGCCA TATGGTGATG CGCGGTTCAC TGTTGGTAAC TAATTAGGA AAAGGGGCGC	6360
TTGTCGGACA TCGCGTGCAA CAGTCGCAGG TGTATCGTAT GTCCCGTGC GATGCGTTAC	6420
CAAAAGTTGC TGCAATCACG TCGAACCGAG TTGTCAGCGT GAATTACGAT GGTTCAAGTT	6480
CGTGGTATGA AGGGCACGGT GCGACATTGA AAGCAACCGA ATTTATCCGG ACCGAAGATT	6540
TTTGAACGGG TACACAAGGT GCGGTGTATT TTGTAATTG GCACGGTGGT ATGAATGCTT	6600
CCTAGTTGGT CTTGACAGGG AGCTCCTTCT CGGGGGAGGA TGGCGGGGT AGATGTTGGT	6660
TCGCTACAGT TACGATGCAA AGGGAAAGGCG GTTGGGGCGT GCGCTGGTGT ACACTGAGTC	6720
GGAGCACCGT ATACCTCGGC AGACCGTTGA CGCTGGGCG ATAAGGgTTG TAGAAGCGCT	6780
GGTGGGTGCG GGGTATGAAA CCTATATCGT CGgTGGGCG gTAAGGGAcc TGGTTGCGGG	6840
AAGGACACCA AAAGATTTG ACATTGTTAC AGGCGCAGTT CCCTCTAGGA TTCTGTAGGTT	6900
GTTCAAGAAC TCGGCCATTA TCGGCAGGCC CTTCCGCATT GTTCATGTGT CGTGTGGCTC	6960
GCAGCTGTAC GAGGTTTCCA CCTTTCGCTC TCGTGTGGGG GAAGggTTCGG TGTGTGTTCC	7020
TGGCACGTTG GAGGAAGATG CATGGCGGAG GGACTTTAGT GTCAATGCCT TGTACTATGA	7080
TCCTCTGAGA AATGTGGTGA TCGATTGTGT CGGTGGAATG GTTGATCTGA AGAGGGCGTCG	7140
CGTGCAGCCG CTCATACCTC TCGGGTCCAT CTTTGTAGAG GACCCAGTGC GCATGCTCCG	7200
GGCATTGAAG TGCTCGGTGA TGTGCGAGTC TTCCATCCCT TTTCTGTCC GCCGCAATT	7260
CGCCGCATGT TTCCCTTCTt GGGGGGTGCT CTCCCTCCCG GTTGACCGAC GAATTtGTAA	7320
AAATCCTCTT TtCCGGTCGG AGCGCcCGC TTGTGCGCGC CCTATGTGGG TAmCAGCTCC	7380
TTCTGTACTT GCAGCCGTCT GTGCACTACT TTATG	7415

(2) INFORMATION FOR SEQ ID NO: 6:

- (i) SEQUENCE CHARACTERISTICS:
- (A) LENGTH: 5271 base pairs
 - (B) TYPE: nucleic acid
 - (C) STRANDEDNESS: double

(D) TOPOLOGY: linear

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 6:

CTTTGATTGG TAAGGTAAGA GAACCGTCAA GAAATAATCC TTTTAAAGTT TTTAGTTCAA	60
AAATAAAAACG CGTTCCACTT ATTTCAATAG TTTTCGTTTC TCCTTCAGTT CTAGAATAGC	120
GAACGTTTCC TGTGGAATAC AGAATTGGG CAGAGGTGTT GTAGACAGTT TGATCTGAAA	180
GAATCGTTGC GGTAACGcTC CCATCATCAA TGGAAATGGA AACATTCCCG GTAAAGACTA	240
CCAACTGATC TTGCAAATCA AAGGGGGACA GCGGCACGCG gCCGGTACTT TCCAATATAG	300
GTTGTCCAGT TTCAGAAAGG CGCGTAGTTT CCTGTGCCGA ATTAATAATA ATTCTTAATT	360
TGCGTAGACC ACTTTCACCA AAGAGGGGAC AAAAAAGAAA GAATATGCC CAAATTGGAT	420
ACCATGCTCT CATGCTTTT TACACAGCAC CTCTTGATG TACCAACCGG TGCGAGACTC	480
AGTTATCTGG GATACTGCTT CAGGGCTTCC CTGTGCAACG ATAGTTCCAC CGTGCATTCC	540
TCCTTCAGGA CCTAAATCGA TAACACAGTC TGCCCTGAACA ATAACATCCA TGTTATGTT	600
GATCATCACA ACCGTATTC CCTGATCTAC CAAGCGTTGA ACAACCTCCA TTAATTGGAT	660
GATATCGGCA AAATGCAATC CCGTAGTAGG TTCTGTCAAAG ATATAGAGAG TTTTTCTGT	720
CGCACGCTTT GAAAGCTCAA GTGCAAGTTT AACGCGCTGG GCTTCTCCCC CTGACAAACGT	780
CAGAGcAGAC TGTCTTAAGC CCACATACCC AAGCCCCACC GAGcAGAGAG CTTCTAGCTT	840
TCGTACTATA GGrGGaACAG CAGAAAAAAA AgAACGSGCT TCTTCGATCG TCATGTCCAG	900
CACATGGGAA ATGTTCTTGc CCTTATAAAA CACAGCTAAT GTCTCCGGT TAAACCGGGT	960
GCCGTGACAC ACATCACAGG TAATGTACAC ATCAGGTAAA AAATTCTATT CAATAGTGAT	1020
AACGCCATCA CCTTTACAAT GCTCACACCG TCCTCCAGGA ACATTGAAAG AAAAACGTCC	1080
TGGTTTATAT CCCCGATTT TTGCTTCAGG AACcTGGGAG AACAGCATTc TAAATATCTGT	1140
AAACACACCC ACATAaGTTG CAGGATTGAG AcGAGGAGTT CTCCCGATAG GACTTTGGTC	1200
TACATAAATT ACTTTATCTA AATGCTCCGT CCCCTCAATC GAGGAAAATT TTCCCTTcAGG	1260
AAAGTCTGCCG TTCATCACAC GGTTGTATAA CGCAGGATAT AGCACATCAA TTAAAAGCGT	1320
TGATTTACCC GAGCCGGATA CTCCGGTAAT GCAGGTAAA GTACCGAGTC GAATACGTAC	1380
AGAAATGTGT TGCAAGTTAT GTTCATGGAC GTCATGCACC GTAAGAACAT TTCCATTTC	1440
CGTTCTCCGT ACTGCAGGAA TGGGTAAATGT AATTGCACCG GCAAGATACT GACCACTAAG	1500
ACTTGCTTGC ACCTGCATAA CTTCAGGTGG ACykCtGCGG CGACAAACATA TCCTCCGTGA	1560

208

ACACCCGCAC CGGGGCCGAG ATCTACAATA TAATCTGCTA CGCGGAGcgT TTgCTCATCG	1620
TGCTCTACCA CAAGCACTGT GTTTCCCAA TCACGCAAAT GAAGAACGCT TTGGATCAGT	1680
CGTTCATTAT CCCGCTGATG CAAACCAATA GACGGTCGT CCAGTATGTA CAAAACCCCT	1740
GTAAGGCCG AACCTATCTG GGTTGCCAGT CTAATTGCTT GTGCTTCTCC GCCGGATAAC	1800
GTGGCAGCAG CCCGTTCAA GGTGAGATAT CCAAGACCCA CGTTCTGAAG AAACCTCTAGG	1860
CGATCGGTAA TTTCTTCAG GATCTGTTGC GCAATTGTCG CTTCTACTTC TGTCAGATGG	1920
AGAGTTTTAA AAAACTCACA CGAACATCATCT ACAGACAGCg CACTGAGTGC GTGGATGTTT	1980
TTTTTTCTA TAGTCACCGC AAGCGACTCT GGCTTTAACG GCATCCCTCG ACACGCTTCA	2040
CATGTACGCA CCGATAAAATA CCGTTCATAT ACCTCGCGCT GTGAGTGAGT ACATGACTCT	2100
GCGTATCTCC TGTGCAGCTC GCTAAAAATT CCCGGCCACG GCTTAATGTA GCGTGCAGGT	2160
CGAGAGCCAT CTTTCGTTTC ATGGGAAAAC TCAAGAGCCT CGCTGCCACT TCCATGCAGG	2220
ATAATATCCA GTGCGTGTGTT TGACAGATTG CGTACCGGAT CATCGAGAGA AAAATGGTAC	2280
TTTTCTGCAG GTGCAGCAAA CCGCACACGG TTCCACTCAT GCTCAGGTTT AAATGGCAAA	2340
AAAGCACCCCT CGTTAAAAGA ACGGTTTTGA TCAGGGACAA TGCGATCTAA ATCAAATGTC	2400
TGCATAATCC CCAGCTCTGC ACAGCTCGGA CAGGCACCAA AAGGTGCGTT AAAAGAGAAC	2460
AAGCGAGGCT GCAATTCGGG TACGGAGACA TTACAGTGCG CGCACCGCTT TTTTGCGAA	2520
AAAATAACT CAGACGGCAG GAGAGCAGAT GTTTCTATCT TTCCGGAAAC GGTCCCAGAA	2580
TTCTCTCCCT GCACTAAGAC GGTCAACAGC CCATCTGCAT AGCCTAGCGT CGTCTCTACT	2640
GATTCTGTTA ATCGTTTACG TACTGTATCT GACAATTGAA TTCTATCGAC AACTATATCG	2700
ATAGAATGCT TTTTTTGCTT ATCCAACGAA ATGCGCTCGT GTAAGTGGAG CAAAGCCCCG	2760
TCAATACGAG CTCGTACAAA ACCATCTTG CGTGCAGCTT CCAAGACCTT GTGGTGTGTA	2820
CTTTTTTTTC CTCGCACCAC CGGGGCAAGC AACTGAATTG TGCTTCCCGA CGGCACGGTC	2880
ATGAGGGTAT CAACAATTG GTCAACGGTT TGTTCTGA TCTCCCGCGC ACAGTGCAGGA	2940
CAATGCCGCG GTCCTATGCG GGCAACAGC AGACGATAGT AGTCATAAAT TTCTGTGACA	3000
GTACCAACCG TTGAGCGAGG GTTACGCTGc GTAGTTTTT GCTCGATGGC AATCGCAGGA	3060
GAAAGACCCCT CGATAGAGTC AACATCCGGC TTATCTAAC GACCTAAAAA CTGGCGAGCG	3120
TATGCAGAAA GGGACTCCAC ATACCGACGC TGTCCCTCTG CAAAAATAGT ATCAAACGCA	3180
AGCGAACTCT TGCCTGAACC AGAAAGACCG GAGATCACCA CAAGCGCATC TCGCGGCAAC	3240
ATAACATCAA TATTCTTCAG ATTATGCTCA CGCGCACCCCT TTATACACAG ATTACGAGCA	3300

209

GCAGAGCCCCA CGCGAACGTC CTGCGACACG TTCTCCCTTT CTACCGGATC CCCATCCATA	3360
AGAGGGCAGA CTATGCGAA TTTTTATGCT TTATGCAATT CCACTCCCTT CGCGGCAGAC	3420
CGCATTACAC CGTCCCTCCA AGCAACGCGA CAATTTTTTC ACTCTGCGCA GACTGTACGA	3480
TGAACATCAG CACACTTCCC TCAAGCAGAA CCGTATCTCC TGAAGGAATG AAAGAGCCAC	3540
GCACCGTTGA GATGAGCAGC ACCAAGAAGC TTCCGTGCAC AGCGATATCC TTCAGGCGCT	3600
TGCCTACCAAG GGGGGACTGC GCGGAAATAG CGAACTCAAC GATTTCTAGC GATCCGTGCG	3660
CGATAGTATG TATGCCAGTG ACGTGGGAAC CGGCCAAATG GCTCATATAA GCGTCAACCA	3720
cTACGTCTTG ATAAGAAAACA GCAACATCGA TGCCAATTTC CCCCCAATA TCCTCCATAA	3780
GGGAGCTGTG TACCAATGCA ACAGCCCGAG CCACTCCGAG CGTCTTCATG TATGCGGCTG	3840
TAATCATATT CAGCTCATAG TTATTAGTGG TGGTAATCAC CAGATCAAAC GTGTCCGGCG	3900
TAATCTCTGC GAAAAAAGCC TCATCTGTGA CATCACCATG ATAGGCAGTA ACGTGCAGAA	3960
ATTGAGCACA CACTGCCTGG GTTGcCtTTC ACTCTTATCA ACCAATACAA GACTCGCACG	4020
CTCCCTTGGA GAAAGACTGA AGGCACTACT GAAAAAGTGC GGCTTGCATT TTTCTGCTAC	4080
ATCCTGTGCC ACGAGCGTAC CTACCGCGCT CATGCCAATG AGTCAATT TTTTACCGG	4140
ATGTATTTTA AAACCGCCA GCTCATAAAA ACGTCCCATG TGTCAGGCG CACAGAGTAC	4200
TGACAGGGCGC ATACCAGAAG CGACCATGGT CTCCCCGTGAG GGAATTACAC TCCTCCCCCG	4260
AACTTCAAAA GCAACGGCAA CAAAAGAAAT TTTTACCAAGA CGACGCATAT CAGAGAGCGT	4320
GATACCATCG AGGCCGCTGC CCTTTGCAAT AGGAAAACGG GCAATTTCAT ACGGTGCATT	4380
TTTCAATGGG ATGACATCGC TGATGGCACC CTGCTCGACG GTGCTCACTA CCGCACCGCAT	4440
CGCTTCCTTA TCCgCAGATA TGAGAAAGTC AATACCAAAA ATACAGCGCG ACTCACGACA	4500
CACCGCGTGA GcGTAGTGGT CATCGTGCCTT TTGGGCTATT TTAATCACTC CGGCATTCAA	4560
GTCGGCGGCT ATACCACACA GTACTATGTT AAGTTCGTCA ACCTCGGTGA CCGCAACAAA	4620
CGCCCTGTGCC TTTGGGATAC CTGCGTCAAC CAGGGTAGCg CGGTGATCTT TTTGATGACG	4680
CaCGAGCGCT TGTGCCCCCG CGGGAAATTTT ACGGGGGACA GGCTCATGCG CGGCAGCAAC	4740
AAGCGTAACC TGATGTCCCC TCGCGCTCAA ACGACCGCGTA AGTTACGCCC CATTGTAACC	4800
GCATCCAACA ACAATAACCC TCATGTCCGA AGGCCATAGT AGCACGAAAT TTTTTGCA	4860
GGCCAGCGCG cAGAACAcCGg CGcACAACGC CTGCCACTCA TATCTTTTTC AAAAGTACCA	4920
CTACCTGTGC GGTAACCGCC GCACCAAGATC CAACAGGTCC GAGGGCTTCG GCAGTCTTG	4980
CCTTAACAAA AACACGTGTT ACGTGCCTGT CCAGGGCCTG CGycaGcGA TGCGCGCATC	5040

210

GCTTCCCGAA ATGGGTGTA TGCAGGCTGC TCAAGACAGA CAACAGCATC GAGATTCA	5100
AGCCGCCAGn CACTGCGCG ACCAGTTGCC AGGTATGGCG GAGCAACGCG CAAGAATGTG	5160
CGTCCTTCCA TCGTCCGTCA CAAGAGGGGA AAAACGTGCC AATATCCCC AGGCCCTTTC	5220
GCCAGCTGGC GTAATAGCGA AGAGGCCCGC ACCGATCGCC CTTCCAAACA G	5271

(2) INFORMATION FOR SEQ ID NO: 7:

- (i) SEQUENCE CHARACTERISTICS:
- (A) LENGTH: 646 base pairs
 - (B) TYPE: nucleic acid
 - (C) STRANDEDNESS: double
 - (D) TOPOLOGY: linear

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 7:

AAGTCCTCT CGTAGCGCTT ACTGTCGGTG AAGCGGGAGT GGAGGACTTG GCCCGGCAGA	60
TGCAGCAGCC CTCTATGCTC CGGTACATCT TTGGCTTGTG GAAACCAAAC ATCTTCGCTT	120
CTGGCCAGGT CTGGAAAGAC AAGGCAACAG ACACAGCTCT GCCGAGGCC TCACAAGAAA	180
CCATCCGGGG AGCCTGTGCT GCTCCACCCC GGGCCCACCC AGGCCGTCA GACAGCAGGA	240
CTCTGTGGAA GGTGGCCTGG ACCCGCCTCC GCTCCTGGCG CTGGCACGGC AAGTGTATGA	300
CACACAGAAC AGCTCAGGTG TTCAGGGAGG CCCCCCGCTC TCAGCACTCC CCCACCCCTG	360
CCCAGCAAAC ATCCTTTCTG AAAATGAGGA AGGGGAGGCT GGTTGGTTTG TTGGCAGGGA	420
GCCAAGCACT TGAGCCATCA TCTGCTGCC CCCCAGGTCC ACGTGAGAAC GAAGCTGGAA	480
TCGGGAGTGG ATCAAGGATT GGAACCCAGG CACTTGGCTA CAGGATATGC TACAAGCTCT	540
CCTGATAATC CTGTAAAATG ATGAAATCAT TTAGGATGTA TCCTGAAATC TGAGACAGG	600
CATACCTTTT CTTCTTGCAT CTTTGAAAGT GaACCCCCC CCACGC	646

(2) INFORMATION FOR SEQ ID NO: 8:

- (i) SEQUENCE CHARACTERISTICS:
- (A) LENGTH: 28295 base pairs
 - (B) TYPE: nucleic acid
 - (C) STRANDEDNESS: double
 - (D) TOPOLOGY: linear

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 8:

GTTCCCTAAA GAAGGAAATG TTTTCTCCtG TGTGTGCAgT CAATGTGCCG GCGTATATyC	60
CGTTAGAATG GGTGCGCTCC AGTGTCAAGG CGATATCCTG ATAGCGCATT CTGAGCGCTC	120

211

GGATATTTT TTTAACACA TTAGTCATAT GCTGAGCTTG TTCCTAAGTT GAACGATTTC	180
ATGGCGCACC TTGTGAATGA GAGTTTGTGT TCCCTGAGAA GGGTTTTTC TCAGGTGAAG	240
TAGCTCAGCG TATGCAAGAA ACTCAAAAGC CTCTTTCT ATGCCGGAGC TGGCTATACT	300
GGATGCGATA CGTGCATGGT GTGsACTGCG TGTGATGCGT TCAAAGTAGT GGCGGAgCAC	360
TGCCTTTTG TGGATGCGCG AAATCGGCGG CCCTTCTGAA AAGGTGAATG TATCCTTTT	420
GGCGAAGGCG TGTTGGGAGA AATCTTCATA CGTGATGGTA CGTACATTGG AAGGTGGTGC	480
GTGTTCTGTA TTTAGCCGTG TAAGTTCTG GATTTCTCT TGTGGGAGAC TGTAAAACCA	540
GTGTGCGTAT GTTTCTAGTG AGCGATTGTC AAAATTCTC GCAATGCAGt CGCAAGTGGC	600
TCAGTTCTCT TGGGATTAAA ACGTTTGAAG GAAGCATGTG GACGTGCGTG TTGGAATCCC	660
TTTCCCCGGT GCAAATCGAG TCCGGCCAAA TAAATAGAAT GTACCCCACA CTGTAAAAGA	720
TATTCGAGTG CAGTCCCCAT GACACTCCC TGTCTTGTG CAGAAACCGA AGGA ^t TGTGCA	780
GGTGTGCGAG AAAGAAGTGT TCTGTATGCG AATGGTAGTT GAGCAAAGAA ATTGGAGAGT	840
ATTTAACACAC ACgtGGAGGA ATGCACGCTT CTAAGGGAAA CAACACCGA AGGGTAGGCG	900
CTGGCGGAAA ATGCTCTGCG GCCCAAAAC TGCCGTCCGT GCTCATGCAG ATATCAGGAG	960
AAATATTGCG GTAGAGAAGT GCCTGCAGTG CTGAGGAGAC CGCAACTATG GGAAGTCCGA	1020
CTCGGGTGAT TTGCTCGAGT CCGAAACCTG cAGCCACACA CAACACTTCC GGTGCGGTCA	1080
GGCACACCGT ACCGGGcGCT CAAGAAAAAA TACATTCTC AGTGTATTCA GTAACCAACG	1140
TTTTCCGAAG TATATGCGCG TGCGATTTC ACTCTGAATT ACTTGATGG TATAgGTGAT	1200
TTCTTGCCAT GTGGACTGAG CTTGTGCGCG CCATATTGG TTTGCTGGCT CCCAGGGAAC	1260
AAACGCGGTT TGCCCCAGCA ATTGCTCAGG GATATGATTA ATTAAGAAAG AATGCAAAC	1320
GCCGCAGTCT GGTCTCCAGA CTGCGTCCCA CTGGATATCG GATGAAGTGA ACGCATCGTT	1380
TGTGTATCGA ACTGCGATAA GCTTTGCATG TGGAAAGCGT GCCCGAAAA ACTCCGCGGT	1440
GTACGACTCA CCCGGCTCTG TTATTACAC AATGCGGCGA CCTTCTAGGA TCGCAGCGAC	1500
AAAACGCTCC GCCTCCCTTC TTGGGTTATA TTTTGAGTGG AGGTTCAGCG ACACGTCTCC	1560
CCGGTTTCGT ACTGCAAGAA ACGGAGGACG TCCGTCTGCG CGCGAGTAGT TCCACTGCTT	1620
GTGATGAGAA GTGAATTCTC CGCGTGGTTA ACAAAAGTACG CGCTGACGCC ACGGGTAAGT	1680
TGGTTTAGGT ACAGCTCAGG ATCGACGGGA AAGGAAGAAA AGACCACAAAC GCGCActTCC	1740
CGGGAAAGCGA CGCATGAGGC AAGGTTGTG TGCCCTACGA GCGTTAATAG CTGGTTGCTC	1800
AGCGCACAGC AGCGGGGTAC CATACTCAGC TCTTCCTGAA GAACGGCGGG GGGCAACGCA	1860

212

CTcTGTTGAGT GAGCTGGTGC CAACTCCGGC TGTcGCCGAT ACACGGGCCT ACTCGCCTCG	1920
AGCTGGGAGC AGCTCGGTGC CAACTCCGGC TGTcGCCGAT ACACGGGCCT ACTCGCCTCG	1980
ATAAACGAGA AGCGCGCCTG TACGGCaAGG ATAAAGGAAC GCACGCGCTC CTCGCCTGAT	2040
GCaTCACC GT GAAGGTCTGC aGCGAGGGCA CCTTCCTCGA CGTCGGTTCG GGTATCAAAG	2100
AGCACTAGGT ATACCGCAGC GCTTCCAGCG GAAAACGGAT ACAAGTGGAG CGCACGCAGC	2160
GTGCTGAATA GATGGGAAGA GAAGTAACCC TGCAATT CGG TAAGTTCCCTG CCTGAACAGG	2220
GTCGTCCACT GCTCCGCAGg CAGGCCGGC GAAAGGAGGG AGGTGGGAAC ACGCATACGG	2280
TAAAACGTGG TGGTGTCAAT CCCGAAGGGa AACCGAGTC AAACATGCAC TCTTCAGTGA	2340
GAAAGAGAAA ACCCGCGCGG GGAATTCCGA GAGATTCAAC CAGCTGTACG AACTGCAGCT	2400
CCAGGCGGTA GTACTCGCAG GAGACACCCG AGCTACCGG GATGCCGGAA GCGCGCGCGA	2460
TAAGACCCAT TAGGAAATCC CTAGCTcGTC AAAGAGGTGT TTGTACGTGT GGAAATACTC	2520
CGATCGCGCA AACTCCTCGA TCTTCTCCTC GGGCAGACTT TCGAGCAGCC GGTCCATATA	2580
CCCGAGCACCA GAGCCACCT CGTCGGCAAG GCTCTTGAG AGGGGAGGTG CAAGGTTGGC	2640
AGAATCCCCC TGCGCAGGAC GGGCCGCTTC CGATTGGGGC GGGGACTCGA TTGCCACCGG	2700
CGTGGCAACC TCCGAAGACT CAGAGCCGGC GGTTCAATC GCGTGGGGAG CAGAAAATC	2760
TAGCGAGGGA ACGGAGACGT CCAACTCCTG CTCTTCTGGC AGTGGGAACG ATGTATCTTC	2820
GCGCTGGGTG TCTTCGTCGT CGAACGACATT TTCAGTATTG AGGGGCTGCT CGGCACCGGC	2880
TGCATCGAGG CCGGTGGAGG AACCGTGCAG CCGTTGAGCC TCAGGCGATG CGTCATCTTT	2940
GGCAACAGGA GACGGATCGG CGCGAGCGTC ACGATCCGG TGAGCCCCCT GCTCTTCAC	3000
ATCGAAAAT TCGTTTCAA GAGTAGCACCG CCCCTGTGTG ACCCACTGCT CCtGcgCGAC	3060
GCgCGCTGcA GCGTGGTCGT CCTCCTCTAG GTAGCTGAAA TCGCCTGCAG CGGACCCGAT	3120
GTGAGAGGAC TCAAAACAC TGAGTTCAAC ACCAAATGCG TCGGCCGCAG AGTCTTTCCC	3180
ATCTTCTCG GTAAACTCAG AGGTGATGAG GATATTGTT AGCTCATCGT TGGTAAGGGC	3240
AATCGTCTCG TCTGGATCAT CGTCACAGAA AAACCCGGAG TAGGCAGCCT CTGTTCTTC	3300
CGGACCCGGA CGGGTGGCGG GAGCCTGAGC AGGCTGGCG CCAGCTTGCC GGGAGAAGGT	3360
TCCCTTCAGC TGGTCTAGAT CAGCGCGGAT ACTGCTGATT TcCTGTGAA TTTTCAGAAG	3420
CAAATCGGTG GAAGCATCGC CTGAGGTACT GGCTGCAGCC CTGCGCGGCC ACCTGGGCAT	3480
CACCCATGAG GGACTGCGCa ATGCGTCCAC GTCACTGAAT GGGGCGCCGC GTACCGGTCC	3540
TTCCGAAGAA GGCTCCGAAG CaGaGTAArt GATAGGGTCG GAATCGACCG GGTATTCCGG	3600

213

TGGCGTCGTGG	CTCAAACTCG	AGAGCAAGTC	GTCGAACCTCT	GTGGTGTGCGC	ACGACCCGCT	3660
GGGGCAGGGC	ATGCTGCACC	CCGCAGGGCTC	TGGCGGCTCT	TCGGTCACAA	CGGAGCTCAT	3720
CGCTGCCGAT	AGATCAACTC	TGTCGTCTG	AGACTGCGCG	GTGCGTGTGA	CCAAATCGTC	3780
GAACGACGGC	ACCTCCGAGT	AAGCGGCATC	TCCGGtGGCG	CATCTGCCGC	CTGGTCGGAG	3840
GAGGACTCGC	CACCCGGAGC	GTGGCTGCAC	CCGCGAGCGA	GCACCTCGTC	CTCCACGGAG	3900
AGGTCAACAA	AGCCACGAGA	ATCCTCTGCA	GAGGACTGCA	CAGGAAGCGC	ACGGTACACA	3960
TGGTGCCGCA	CCGCATCCAC	GCGCGCACCC	GGGGGAGCGT	GTTGGGCAGG	CGCCTCTGGA	4020
TCCCTCTGGGC	CACCTGAGCT	GACGGGTGGC	TCCTGGCCTT	CCGGCGGTAC	CTTTACCCAC	4080
ACGCCACACG	CATCCCAGGC	ACGGTCTTCC	CCCTGCTCCT	GAGGGCGTGC	AGTTGCGTGT	4140
GGACTGTCCA	TTTCTGCTGA	CTCTGATCCC	ATAGCATCTT	CCTCTGTGCC	GTGAACTCCC	4200
TCACACATCC	TGTGTCGGCA	GCCGTGCCGC	AAAGCATGAG	CATAGCGTGT	CGGGTCCGTT	4260
TTTTCAACAA	TTCGTAAAAG	CACCCCGTTC	TTACGCCGTA	AATGACGGCA	ACGGCGGGGT	4320
GCGGAAGGCA	CCTGCAGTGC	CTGTGTCATT	ACCTCCTGGC	GGGGGGGGGT	ACGGACAAAAA	4380
ACAGGGGTGA	AGATGTGGAA	GACAGTGCAG	ACACGCAGTC	AGAGGAAGCG	GTGAGGACTC	4440
TATGCGCATT	TACTTGAGGG	TAGTACTTCC	CCTGTCTCTT	GCGCTGAACA	GCTACGGTGT	4500
ACTCGCCTTT	TTCTGGGGAG	AGCGGGGGGT	GTGTGCCATG	CGGCTACTGG	AACGTGAGAA	4560
AAAGGAGCTC	GTCCATCACA	TCCAGACGCT	CGCAGAGCGT	GGGCGCGACT	TGGCTGCGGT	4620
GGTGGACGCT	CTATCCTTTG	ACGAAGAGAC	TATCGGTGCG	TATGCGCGTc	AGCTAGGATA	4680
TGTCCGCGCG	GGGGATGTGT	TAGTGAGGCC	GGTAAACTTT	ACCTTTGCGC	ACATGCATAC	4740
CCTgATTCTG	GGGATGCACG	TCCGCTtGTT	GCACCTGCAT	GTTTTAGCGA	CACGCGTTGA	4800
CAAGGTGTAC	GCGCTGTkcs	TcGGCTTcTT	TgTCGTCTG	TTACAGCTGC	TGTGGGGTAG	4860
CCCGCGTGC	TATTTAAAAA	CATGAGGCAGC	AGTCTGCAGC	CTnTGCgCgT	GCGCTCAAGG	4920
CCGGTGCCT	CGTAGCGTTG	CCGACAGATA	CGGTGTACGG	TTTCTCTGGC	CTTGTGCCAC	4980
ACGCTGTTCC	GGATCTCATA	TGTCTGAAGG	CGCGTGGGTG	CACAGAGACG	GAAGGGAAC	5040
GGAGAGAGGG	CTATCCGTTTC	ATTGCACTGC	TTGCAGATCC	ACAGGACGTG	GTTGTCTATA	5100
CCGGGACGCG	GCTTCTGCGG	AGTTTCGTGC	GCTGTGGCCT	GGCCCGTATA	CGTTTGtntG	5160
CGCATGCAAG	ACGGCGCGAC	GCAGGcgTTC	CGCTGTCTG	CTGACCTGTG	cTGCCTCAG	5220
TGATACGGGC	AGTTGGGGGA	GGCATCTTTT	CCACGAGTGC	AAATCGGCAC	GGCGAGCCGC	5280
CGCTGCAAGA	TGCACAGGAC	ATCGACCACA	TCTTTGGAAA	GCATCTTGCG	CTGACCGTAG	5340

214

ACGCAGGACC ACTGACCGGC TCCCCAAGCG CGGTGATAGA CCTCACGCAC CCCGTGcCGC	5400
GTGTGctCCG CGCTGGTGC GCGCCGTTGC CTCTTGCAAG ACTGGAAAGG CGTGACTCTC	5460
CTTCCCTCCC TCATGTGGGG GAAGTATGTA AAGAATGAGG TCAGTTGCGG ACCCACTCCT	5520
CAGGCACCTC TTCATACTCA ACAGAACGCG TGTGCCCTGC GCCCGCACG ACAGGGGCGG	5580
TTCTGACGCC TGAAAGCCGG TTGTTTCCTG AGATCTTGG GAACCTCCAC TGGAGCGGAC	5640
GCGCCTCcTG CGCGAGCTTC TGAGCGATCT CAGGGGGAAA ATGAGAAAAC TGTCGCGCAT	5700
TGACTGCCCT CTTCTGTGAC ACCACGAGCA CTCCGCTGA AATTTTCAGC TTGAGGGAA	5760
CGGAAAGTCC GGAGCGCAAC ACCGCGATGC CGCGGCCCTTC GCTCACGATG ACGATTCTTT	5820
CCCCCTCCTC TGTGCTGTGC CAGGAACCGG CGAGGGCATC TAGGGACGAA ACGGATTCTT	5880
CACCTGCGCG CGTGTGCCGC ATGCTAGACG TCTCTGTCTG ATTTCTGTTC AGCGGGACAG	5940
AACGGTCAAA CACATCTCGG ACCAGGTGGC GCGAGTCAG CAGAACGCG CCGCCGTCT	6000
CATATGTTT GGAGAGCAGC CGCGTGGCGT TGTGATCTTT AmCCTTGAGT GCAACAGCCA	6060
GTCTAAATCCC CTCAGGGGTG ArGTCCATCG CGCCGCAAAA TATGTAGTCAG ACGATTCCCTT	6120
TTTCGGGAAA ACGGTGCCGC ACCGCTTGCT CTCTGCAATC TACAAmGTGA TACCCACGCA	6180
ACTCCCGAAT GAAgAAAAG ACCGCGTCGT TGATGGTAGT TTCTGTGTGC GCAGGGCACAC	6240
CAgACACTTc TAGCCTGTAG ACCCAAACAC GAGGAGCTGC GTGAACAGCA TCGCGGAACA	6300
GCACGAACAG GATAAGCAGC CGAGAAGAAC GGACACCTTT TTTCATGAGA CTAGTGGTGT	6360
CGCTCACAGA GGCTGCCGG ACGCTCCCGT GCGTTGTCGC GAGCTTGTAT CTGCGCGCGC	6420
TTGTCAAAAA GCTTCTTGCC CTTCCAGATT CCCAGCGCTA CCTTCACCCG CCCTGCTTTT	6480
AGGTAAAACCT CCAGGGGGAC CAGAGTATAG CCTTTCTCTT CAACCTTGCG CTTCAAGCGC	6540
GCAATCTGGT CCCGATGTGC CAGTAACCTTC CGCATCCGAT CCGGATTGGG GGCAAAGGAG	6600
CAAGCATGCA CGTACTCCGC AATATGCACA TTCTTTAGCC ACAGCTGCC TCCGCGCATC	6660
TCTGcAAATG CGTCAGGAAA AGAAAGATGC CCCGCGGCCA CAGACTTCAC CTCCGTGCCT	6720
TCAAGCGCGA TGCCACACTC TAGACGGTCT TCCACATGGT AATTGAAAAA AGCcTTGCGG	6780
TTCTTTGCAA TGAGATGGGT TCCTGTGCC CTCATGGCGC CGGATGCTAC CGGATAGGCA	6840
CTTCCCTTGT CAATTGATT ATCGCCGTGT TAGGCTGCCG TGTCTGGAG GGACGCCGTT	6900
TTATGTTGC GCGGTGGAGA AGGTACTCAT ATTTGGCGCG GCGCGAAGCA CGGCAGGAATG	6960
CGACCCGAGT TTGTAGTgCT GGGGTGGGCT TCTTTCTGTT CTATCTTTT ATCACTACGC	7020
ATGTGGTTGC AGCGTATCGC ATTCAAGGCGG ACTCGATGCA GCCGACCCCTG AGCGCAGGGG	7080

215

ATTGCGTTCT	TGCCCTCGTCC	CTGTTTCGCT	TTGCCCGCAT	CAAGCGGGGG	GATTTGGTGC	7140
TTGCAACTCC	CCTTGAGAAA	GAGGATATAG	GCCTGTAA	AAGGGCGATG	AATGCTGTGT	7200
TnAGgnTCG	CAAGCCTTCA	ATTGTACCGG	CCGTTGGCG	CGGCAGATCG	CATGTTTCG	7260
CGGCCGCAA	TGCGCAGGGT	GGTGGGCCTT	CCAGGGACA	CTGCTATAT	GCGCGATTTT	7320
GTGCTGTACG	TTAAGCCCCA	CGGTCAAGCAA	CACTTCCTCA	CGGAATTGTA	AGTGAGTGCA	7380
GTTAGCTACG	ACGTGCGTAA	GGGGGTGCTT	CCTGAGCATT	GGTCTGAACG	GCTTCCCTTT	7440
TCTGGTTTCA	TGGAAGAGAT	GCAGTTGGAC	GAGCACTCCT	ATTGCTGCTG	TGCGATAATC	7500
GAATTGCTC	CAGTGATTCT	CGTCCTGGG	GTGCCATCGA	CGGTAGTACG	CAGATAAAAG	7560
CAAAGGCATT	CATGCGTTAT	TTCCCTTTCG	GAGCATTGG	TGTCTTGTAG	TGTGTAGGCG	7620
CCGCATTG	GGTGCCTGtG	CCGATCGTGC	TGTTCCCTTT	ATCATGTCCTT	CTGAGGTCCG	7680
TGCGTCTTGT	TACGTGCACA	TCCCCCTCTG	TGCGCAACGC	TGTGCTTACT	GCGATTTTTA	7740
CTCCCTGGTG	CGTTCAACCT	ATTTTAGGCC	TCATCAGCCT	TGTCCGCATT	TTATCGATCG	7800
GCTGCTACAG	GATGTGGCAT	TGCAGCGGGA	GTGCTTGGG	GTCCACGGkT	GGCAGACAGT	7860
GTATATGGGT	GGAGgTACCC	CTTCGCTATT	GGCACCGCAG	GACATTGTC	ATTTTTGCGT	7920
AGCGTTACGC	GCCGCGCAGC	GGTATCCGAT	TCAGGAGTTC	ACTCTTGAGG	TGAATCCTGA	7980
GGATGTGACC	GAAGAGTTT	TGTGTGCGTG	TGCAAGAACG	GGAGTAAACC	GTTTATCCCT	8040
TGGGGTACAA	AGTCTGCGTG	ATGAGGTGTT	CGGTGCGGAG	CGTCGTGCAG	CCTCTGCTGA	8100
ATGTGCTCGT	ACCCGCTCCG	CGTGATGACG	GCAAATGCGC	GCTTTTCTC	TGGCGGGGTG	8160
CGTATTCAG	CAGATCTCAT	CGCTGGATTG	CGCGGGAAA	CGGCCGAAT	GGTGCCTGAG	8220
GATaTAGATG	AGCTTTGTC	TTTTGGCTG	AGACACGTGT	CGCTATATGG	TGTTGTGTGA	8280
CCGCATCCGA	CTGAAACGCA	AGAGGAGCGA	ATTGCAAGCG	TTTGGGCACA	CGGCAGCGCG	8340
TATCTGGTGC	GTGCaGGATT	TAACCGGTAT	GAGCTTCGA	ATTTGCACG	TACTGCgGCG	8400
GACGAGAGCG	CGCACAAACAG	AGCATATTGG	CGGATGGCAC	CGCACGCAGG	GGTGGGGCCT	8460
GGCGCAGTTG	GCACCGTTT	TGTCAACCTT	TCTTTATCAA	AGGAGGGGGC	GTGGCGATC	8520
CCGACACCG	TGCGGAAACA	TCTTGGCCAA	TACTTAGCAG	AAGTGTGTG	GGAAAATGTG	8580
TATGAGCACG	AATTCTTAC	AGAACATATG	TGTGTGCAAG	AAGCATTGTT	AATGGGATTA	8640
CGTCTTGAAC	AGGGACTGGA	TGTGGTTACA	TTTCGTGCGC	GGTCGGGAA	GGGAATTGAA	8700
CGGTACATTG	GCAAAACAAT	CGCGCGGTGG	CAGTGTCAATG	GCCGAATGCA	CGGCACGGCG	8760
ACGTCATTGC	GTTTGAGTGC	GCAGgCACGG	GTATTTCTGG	ACAGTTTTT	GCGAGAGGCC	8820